

**NSWC TR 90-21**

# **NSWC LIBRARY OF MATHEMATICS SUBROUTINES**

**ALFRED H. MORRIS, JR.  
ENGINEERING AND INFORMATION SYSTEMS DEPARTMENT**

**JANUARY 1990**

Approved for public release; distribution unlimited.

DESTRUCTION NOTICE —For classified documents, follow procedures as outlined in Chapter 17 of OPNAVINST 5510.1H. For unclassified, limited documents, destroy by any method that will prevent disclosure of contents or reconstruction of the document.

**NAVAL SURFACE WARFARE CENTER**  
Dahlgren, Virginia 22448-5000 • Silver Spring, Maryland 20903-5000

## FOREWORD

In 1976 development of the NSWC library of general purpose numerical mathematics subroutines began. The subroutines are written in ANSI standard Fortran. This manual describes the subroutines in the 1990 edition of the library. The manual supersedes NSWC TR 86-251 (1987). The development of the NSWC library is funded by the Computing Systems and Networks Division, Engineering and Information Systems Department, NSWC.

Approved by:



R. T. RYLAND, JR., Head  
Engineering and Information  
Systems Department



2000135423

3/1/90

D *seen*

NSWC TR 90-21

D1, D2, D21, D22,  
D23, ~~D24~~, D25, D4  
*rec'd cyp*

D3 *LP*D26 *hmc*D27 *SCP*D28 *✓*

# NSWC LIBRARY OF MATHEMATICS SUBROUTINES *E231 Retae*

BY ALFRED H. MORRIS, JR.

ENGINEERING AND INFORMATION SYSTEMS DEPARTMENT

JANUARY 1990

Approved for public release; distribution unlimited.

DESTRUCTION NOTICE - For classified documents, follow procedures as outlined in Chapter 17 of OPNAVINST 5510.1H. For unclassified, limited documents, destroy by any method that will prevent disclosure of contents or reconstruction of the document.



## NAVAL SURFACE WARFARE CENTER

Dahlgren, Virginia 22448-5000 • Silver Spring, Maryland 20903-5000

NSWC TR 90-21

# CONTENTS

	Page
Introduction .....	1
<b>Elementary Operations</b>	
Machine Constants - SPMPAR,DPMPAR,IPMPAR .....	3
Sorting Lists - ISHELL,SHELL, AORD, RISORT,SHELL2,DSORT, DAORD,DISORT,DDSORT .....	5
Cube Root - CBRT,DCBRT .....	7
Four Quadrant Arctangent - ARTNQ,DARTNQ .....	7
Length of a Two-Dimensional Vector - CPABS,DCPABS .....	7
Reciprocal of a Complex Number - CREC,DCREC .....	9
Square Root of a Double Precision Complex Number - DCSQRT .....	9
Conversion of Polar to Cartesian Coordinates - POCA .....	11
Conversion of Cartesian to Polar Coordinates - CAPO .....	11
Rotation of Axes - ROTA .....	11
Planar Givens Rotations - SROTG,DROTG .....	13
Three Dimensional Rotations - ROT3 .....	15
Rotation of a Point on the Unit Sphere to the North Pole - CONSTR .....	17
Hyperbolic Sine and Cosine Functions - SNHCSH .....	19
Exponentials - REXP,DREXP .....	21
Logarithms - ALNREL,RLOG,RLOG1,DLNREL,DRLOG,DRLOG1 .....	23
<b>Geometry</b>	
Determining if a Point is Inside or Outside a Polygon - LOCPT .....	25
The Convex Hull for a Finite Planar Set - HULL .....	27
Areas of Planar Polygons - PAREA .....	29
Hamiltonian Circuits - HC .....	31
<b>Special Functions</b>	
Error Function - CERF,CERFC,ERF,ERFC,ERFC1,DCERF,DCERFC, DERF,DERFC,DERFC1 .....	35
Inverse Error Function - ERFINV .....	41
Normal Probability Distribution Function - PNDF .....	43
Inverse Normal Probability Distribution Function - PNINV .....	45
Dawson's Integral - DAW .....	47
Complex Fresnel Integral - CFRNL1 .....	49
Real Fresnel Integrals - FRNL .....	51
Exponential Integral Function - CEXPLI,EXPLI,DEI,DEI1 .....	53
Sine and Cosine Integral Functions - SI,CIN .....	57
Dilogarithm Function - CLI,ALI .....	59
Gamma Function - CGAMMA,GAMMA,GAMLN,DCGAMA, DGAMMA,DGAMLN .....	61
Digamma Function - CPSI,PSI,DCPSI,DPSI .....	65



Logarithm of the Beta Function – BETALN,DBETLN .....	67
Incomplete Gamma Ratio Functions – GRATIO,RCOMP .....	69
Inverse Incomplete Gamma Ratio Function – GAMINV .....	71
Incomplete Beta Function – BRATIO,ISUBX,BRCOMP .....	73
Bessel Function $J_\nu(z)$ – CBSSLJ,BSSLJ,BESJ .....	75
Bessel Function $Y_\nu(z)$ – BSSLY .....	77
Modified Bessel Function $I_\nu(z)$ – BSSLI,BESI .....	79
Modified Bessel Function $K_\nu(z)$ – CBSSLK,BSSLK .....	81
Airy functions – CAI,CBI,AI,AIE,BI,BIE .....	83
Complete Complex Elliptic Integrals of the First and Second Kinds – CK,CKE .....	87
Real Elliptic Integrals of the First and Second Kinds – ELLPI,RFVAL,RDVAL,DELLPI,DRFVAL,DRDVAL .....	91
Real Elliptic Integrals of the Third Kind – EPI,RJVAL, DEPI,DRJVAL .....	95
Jacobian Elliptic Functions – ELLPF,ELPFC1 .....	99
Weierstrass Elliptic Function for the Equianharmonic and Lemniscatic Cases – PEQ,PEQ1,PLEM,PLEM1 .....	103
Integral of the Bivariate Density Function over Arbitrary Polygons and Semi-infinite Angular Regions – VALR2 .....	107
Circular Coverage Function – CIRCVC .....	109
Elliptical Coverage Function – PKILL,PKILL3 .....	111

## Polynomials

Copying Polynomials – PLCOPY,DPCOPY .....	113
Addition of Polynomials – PADD,DPADD .....	115
Subtraction of Polynomials – PSUBT,DPSUBST .....	117
Multiplication of Polynomials – PMULT,DPMULT .....	119
Division of Polynomials – PDIV,DPDIV .....	121
Real Powers of Polynomials – PLPWR,DPLPWR .....	123
Inverses of Power Series – PINV,DPINV .....	125
Derivatives and Integrals of Polynomials – MPLNMV .....	127
Evaluation of Chebyshev Expansions – CSEVL,DCSEVL .....	129
Lagrange Polynomials – LGRNGN,LGRNGV,LGRNGX .....	131
Orthogonal Polynomials on Finite Sets – ORTHOS,ORTHOV, ORTHOX .....	133

## Solutions of Nonlinear Equations

Zeros of Continuous Functions – ZEROIN .....	135
Solution of Systems of Nonlinear Equations – HBRD .....	137
Solutions of Quadratic, Cubic, and Quartic Equations – QDCRT,CBCRT,QTCRT,DQDCRT,DCBCRT,DQTCRT .....	139
Double Precision Roots of Polynomials – DRPOLY,DCPOLY .....	141
Accuracy of the Roots of a Real Polynomial – RBND .....	143

## Vectors

Copying Vectors – SCOPY,DCOPY,CCOPY .....	145
-------------------------------------------	-----

Interchanging Vectors – SSWAP,DSWAP,CSWAP .....	147
Planar Rotation of Vectors – SROT,DROT,CSROT .....	149
Dot Products of Vectors – SDOT,DDOT,CDOTC,CDOTU .....	151
Scaling Vectors – SSCAL,DSCAL,CSCAL,CSSCAL .....	153
Vector Addition – SAXPY,DAXPY,CAXPY .....	155
$L_1$ Norm of a Vector – SASUM,DASUM,SCASUM .....	157
$L_2$ Norm of a Vector – SNRM2,DNRM2,SCNRM2 .....	159
$L_\infty$ Norm of a Vector – ISAMAX,IDAMAX,ICAMAX .....	161

## Matrices

Packing and Unpacking Symmetric Matrices – MCVFS,DMCVFS, MCVSF,DMCVSF .....	163
Conversion of Real Matrices to and from Double Precision Form – MCVRD,MCVDR .....	165
Storage of Real Matrices in the Complex Matrix Format – MCVRC .....	167
The Real and Imaginary Parts of a Complex Matrix – CMREAL,CMIMAG .....	169
Copying Matrices – MCOPY,SCOPY,DMCOPY,CMCOPY .....	171
Computation of the Conjugate of a Complex Matrix – CMCONJ .....	173
Transposing Matrices – TPOSE,DTPOSE,CTPOSE,TIP,DTIP,CTIP .....	175
Computing Adjoints of Complex Matrices – CMADJ,CTRANS .....	177
Matrix Addition – MADD,SMADD,DMADD,CMADD .....	179
Matrix Subtraction – MSUBT,MSUBT,DMSUBT,CMSUBT .....	181
Matrix Multiplication – MTMS,DMTMS,CMTMS,MPROD,DMPROD, CMPROD .....	183
Product of a Packed Symmetric Matrix and a Vector – SVPRD,DSVPRD .....	185
Transpose Matrix Products – TMPROD .....	187
Symmetric Matrix Products – SMPROD .....	189
Kronecker Product of Matrices – KPROD,DKPROD,CKPROD .....	191
Inverting General Real Matrices and Solving General Systems of Real Linear Equations – CROUT,KROUT, NPIVOT,MSLV,DMSLV .....	193
Solutions of Real Equations with Iterative Improvement – SLVMP .....	197
Solution of Almost Block Diagonal Systems of Linear Equations – ARCECO,ARCESL .....	199
Solution of Almost Block Tridiagonal Systems of Linear Equations – BTSLV .....	201
Inverting Symmetric Real Matrices and Solving Symmetric Systems of Real Linear Equations – SMSLV,DSMSLV .....	203
Inverting Positive Definite Symmetric Matrices and Solving Positive Definite Symmetric Systems of Linear Equations – PCHOL,DPCHOL .....	207
Solution of Toeplitz Systems of Linear Equations – TOPLX,DTOPLX .....	209

Inverting General Complex Matrices and Solving General Systems of Complex Linear Equations – CMSLV,CMSLV1, DCMSLV .....	211
Solution of Complex Equations with Iterative Improvement – CSLVMP .....	215
Singular Value Decomposition of a Matrix – SSVDC,DSVDC, CSVDC .....	217
Evaluation of the Characteristic Polynomial of a Matrix – DET,DPDET,CDET .....	219
Solution of the Matrix Equation $AX + XB = C$ – ABSLV, DABSLV .....	221
Solution of the Matrix Equation $A^tX + XB = C$ when $C$ is Symmetric – TASLV,DTASLV .....	223
Solution of the Matrix Equation $AX^2 + XB + C = 0$ – SQUINT .....	225
Exponential of a Real Matrix – MEXP,DMEXP .....	227
<b>Large Dense Systems of Linear Equations</b>	
Solving systems of 200–400 Linear Equations – LE,DPLE,CLE .....	229
<b>Banded Matrices</b>	
Band Matrix Storage .....	231
Conversion of Banded Matrices to and from the Standard Format – CVBR,CVBC,CVRB,CVCB,CVRB1,CVCB1 .....	233
Conversion of Banded Matrices to and from Sparse Form – MCVBS,CMCVBS,MCVSB,CMCVSB .....	235
Transposing Banded Matrices – BPOSE,CBPOSE .....	237
Addition of Banded Matrices – BADD,CBADD .....	239
Subtraction of Banded Matrices – BSUBT,CBSUBT .....	241
Multiplication of Banded Matrices – BPROD,CPROD .....	243
Product of a Real Banded Matrix and Vector – BVPRD,BVPRD1, BTPRD,BTPRD1 .....	245
Product of a Complex Banded Matrix and Vector – CBVPD, CBVPD1,CBTPD,CBTPD1 .....	247
Solution of Banded Systems of Real Linear Equations – BSLV,BSLV1 .....	249
Solution of Banded Systems of Complex Linear Equations – CBSLV,CBSLV1 .....	251
<b>Sparse Matrices</b>	
Storage of Sparse Matrices .....	253
Conversion of Sparse Matrices to and from the Standard Format – CVRS,CVDS,CVCS,CVSR,CVSD,CVSC .....	255
Conversion of Sparse Real Matrices to and from Double Precision Form – SCVRD,SCVDR .....	257
The Real and Imaginary Parts of a Sparse Complex Matrix – CSREAL,CSIMAG .....	259

Computing $A + iB$ for Sparse Real Matrices $A$ and $B$ – SCVRC .....	261
Copying Sparse Matrices – RSCOPY, DSCOPY, CSCOPY .....	263
Computing Conjugates of Sparse Complex Matrices – SCONJ .....	265
Transposing Sparse Real Matrices – RPOSE, RPOSE1 .....	267
Transposing Sparse Double Precision Matrices – DPOSE, DPOSE1 .....	269
Transposing Sparse Complex Matrices – CPOSE, CPOSE1 .....	271
Addition of Sparse Matrices – SADD, DSADD, CSADD .....	273
Subtraction of Sparse Matrices – SSUBT, DSSUBT, CSSUBT .....	275
Multiplication of Sparse Matrices – SPROD, DSPROD, CSPROD .....	277
Product of a Real Sparse Matrix and Vector – MVPRD, MVPRD1, MTPRD, MTPRD1 .....	279
Product of a Double Precision Sparse Matrix and Vector – DVPRD, DVPRD1, DTPRD, DTPRD1 .....	281
Product of a Complex Sparse Matrix and Vector – CVPRD, CVPRD1, CTPRD, CTPRD1 .....	283
Ordering the Rows of a Sparse Matrix by Increasing Length – SPORD .....	285
Reordering Sparse Matrices into Block Triangular Form – BLKORD .....	287
Solution of Sparse Systems of Real Linear Equations – SPSLV, RSLV, TSLV .....	289
Double Precision Solution of Sparse Systems of Real Linear Equations – DSPSLV, DSLV, DTSLV .....	293
Solution of Sparse Systems of Complex Linear Equations – CSPSLV, CSLV, CTSLV .....	297

## **Eigenvalues and Eigenvectors**

Computation of Eigenvalues of General Real Matrices – EIG, EIG1 .....	301
Computation of Eigenvalues and Eigenvectors of General Real Matrices – EIGV, EIGV1 .....	303
Double Precision Computation of Eigenvalues of Real Matrices – DEIG .....	305
Double Precision Computation of Eigenvalues and Eigenvectors of Real Matrices – DEIGV .....	307
Computation of Eigenvalues of Symmetric Real Matrices – SEIG, SEIG1 .....	309
Computation of Eigenvalues and Eigenvectors of Symmetric Real Matrices – SEIGV, SEIGV1 .....	311
Computation of Eigenvalues of Complex Matrices – CEIG .....	313
Computation of Eigenvalues and Eigenvectors of Complex Matrices – CEIGV .....	315
Double Precision Computation of Eigenvalues of Complex Matrices – DCEIG .....	317
Double Precision Computation of Eigenvalues and Eigenvectors of Complex Matrices – DCEIGV .....	319

## **$\ell_1$ Solution of Linear Equations**

$\ell_1$ Solution of Systems of Linear Equations with Equality and Inequality Constraints – CL1 .....	321
----------------------------------------------------------------------------------------------------------	-----

## **Least Squares Solution of Linear Equations**

Least Squares Solution of Systems of Linear Equations – LLSQ,HFTI,HFTI2 .....	323
Least Squares Solution of Overdetermined Systems of Linear Equations with Iterative Improvement – LLSQMP .....	327
Double Precision Least Squares Solution of Systems of Linear Equations – DLLSQ,DHFTI,DHFTI2 .....	329
Least Squares Solution of Systems of Linear Equations with Equality and Inequality Constraints – LSEI .....	333
Least Squares Solution of Systems of Linear Equations with Equality and Nonnegativity Constraints – WNNLS .....	337
Least Squares Iterative Improvement Solution of Systems of Linear Equations with Equality Constraints – L2SLV .....	341
Iterative Least Squares Solution of Banded Linear Equations – BLSQ .....	345
Iterative Least Squares Solution of Sparse Linear Equations – SPLSQ,STLSQ .....	347

## **Optimization**

Minimization of Functions of a Single Variable – FMIN .....	349
Minimization of Functions of $n$ Variables – OPTF .....	351
Unconstrained Minimum of the Sum of Squares of Nonlinear Functions – LMDIFF .....	353
Linear Programming – SMPLX,SSPLX .....	355
The Assignment Problem – ASSGN .....	359
0-1 Knapsack Problem – MKP .....	361

## **Transforms**

Inversion of the Laplace Transform – LAINV .....	363
Fast Fourier Transform – FFT,FFT1 .....	367
Multivariate Fast Fourier Transform – MFFT,MFFT1 .....	369
Discrete Cosine and Sine Transforms – COSQI,COSQB,COSQF, SINQB,SINQF .....	371

## **Approximation of Functions**

Rational Minimax Approximation of Functions – CHEBY .....	375
$L_p$ Approximation of Functions – ADAPT .....	377
Calculation of the Taylor Series of a Complex Analytic Function – CPSC,DCPSC .....	381

## **Curve Fitting**

Linear Interpolation – TRP .....	385
Lagrange Interpolation – LTRP .....	387

Hermite Interpolation – HTRP .....	389
Conversion of Real Polynomials from Newton to Taylor Series Form – PCOEFF .....	391
Least Squares Polynomial Fit – PFIT .....	393
Weighted Least Squares Polynomial Fit – WPFIT .....	395
Cubic Spline Interpolation – CBSPL,SPLIFT .....	397
Weighted Least Squares Cubic Spline Fitting – SPFIT .....	399
Cubic Spline Evaluation – SCOMP,SCOMP1,SCOMP2 .....	401
Cubic Spline Evaluation and Differentiation – SEVAL, SEVAL1,SEVAL2 .....	403
Integrals of Cubic Splines – CSINT,CSINT1,CSINT2 .....	405
<i>N</i> -Dimensional Cubic Spline Closed Curve Fitting – CSLOOP, LOPCMP,LOPDF .....	407
Spline under Tension Interpolation – CURV1 .....	409
Spline under Tension Evaluation – CURV2 .....	411
Differentiation and Integration of Splines under Tension – CURVD,CURVI .....	413
Two Dimensional Spline under Tension Curve Fitting – KURV1,KURV2 .....	415
Two Dimensional Spline under Tension Closed Curve Fitting – KURVP1,KURVP2 .....	417
Three Dimensional Spline under Tension Curve Fitting – QURV1,QURV2 .....	419
<i>B</i> -Splines .....	421
Piecewise Polynomial Interpolation – BSTRP .....	423
Conversion of Piecewise Polynomials from <i>B</i> -Spline to Taylor Series Form – BSPP .....	425
Piecewise Polynomial Evaluation – PPVAL .....	427
Weighted Least Squares Piecewise Polynomial Fitting – BSL2 .....	429
<b>Surface Fitting over Rectangular Grids</b>	
Bi-Splines under Tension .....	431
Bi-Spline under Tension Surface Interpolation – SURF .....	433
Bi-Spline under Tension Evaluation – SURF2,NSURF2 .....	435
<b>Surface Fitting over Arbitrarily Positioned Data Points</b>	
Surface Interpolation for Arbitrarily Positioned Data Points – BVIP,BVIP2 .....	437
<b>Manifold Fitting</b>	
Weighted Least Squares Fitting with Polynomials of <i>n</i> Variables – MFIT,DMFIT,MEVAL,DMEVAL .....	441
<b>Numerical Integration</b>	
Evaluation of Integrals over Finite Intervals – QAGS QSUBA,DQAGS .....	445

Evaluation of Integrals over Infinite Intervals – QAGI, DQAGI .....	449
Evaluation of Double Integrals over Triangles – CUBTRI .....	453
<b>Integral Equations</b>	
Solution of Fredholm Integral Equations of the Second Kind – IESLV .....	455
<b>Ordinary Differential Equations/Initial Value Problems</b>	
The Initial Value Solvers – Introductory Comments .....	459
Adaptive Adams Solution of Nonstiff Differential Equations – ODE .....	461
Adaptive RKF Solution of Nonstiff Differential Equations – RKF45 .....	465
Adaptive RKF Solution of Nonstiff Differential Equations with Global Error Estimation – GERK .....	469
Adaptive Solution of Stiff Differential Equations – SFODE, SFODE1 .....	473
Fourth-Order Runge-Kutta – RK .....	477
Eighth-Order Runge-Kutta – RK8 .....	479
<b>Partial Differential Equations</b>	
Separable Second-Order Elliptic Equations on Rectangular Domains – SEPDE .....	481
<b>Random Number Generation</b>	
Uniform Random Number Generator – URNG .....	485
Gaussian Random Number Generator using the Box-Muller Transformation – NRNG .....	487
Appendix. Installation of the NSWC Library .....	489
Index .....	491
Distribution	

## INTRODUCTION

In 1976 development of the NSW library of numerical mathematics subroutines began. The objective was to form a library of general purpose subroutines that would provide a basic computational ability in a variety of mathematical activities. The subroutines were to be written in Fortran. Even though the subroutines were intended for use on the CDC 6000-7000 series computers, emphasis was to be placed on their transportability. Currently, the library is used on a variety of computers, ranging from mainframes such as the CDC Cyber 995 to personal computers, such as the IBM PC. The 1990 edition of the NSW library contains 870+ functions and subroutines. This manual describes the 453 functions and subroutines in the library that are available for general use. (The remaining functions and subroutines are supportive, normally being of little interest to most users.) A brief appendix is included, which provides information for installing the library.

All functions and subroutines are examined before being accepted for the NSW library. Primary considerations are the reliability and transportability of the function or subroutine, its efficiency and ease of use, and its generality. In regard to reliability, the major concerns are accuracy, the stability and robustness of the algorithm being used, and the overall quality of the code. All routines are tested before being accepted for the library. The functions and subroutines in the library are always subject to reexamination and possible modification. When better routines are obtained, the older routines are normally eliminated.

In regard to transportability, it is clear that machine dependent constants and precision dependent algorithms cannot be avoided. However, machine dependent code and I/O statements are not permitted. For a function or subroutine to be acceptable, it is required that the coding satisfy the 1966 and 1977 ANSI Fortran standards, the only exception being the declaration of assumed size arrays, where the standards conflict. In this case, statements such as `REAL A(1)` and `REAL A(*)` are equally acceptable for declaring as arrays arguments *A* of a function or subroutine. Two versions of the code for the NSW library are maintained, providing the appropriate assumed size array declarations for the 1966 and 1977 standards.

The ease of use criterion is of considerable importance. The main purpose of the library is to provide a service to as broad an audience as is possible. Hence, it is important that duplicate abilities be kept to a minimum, and that the functions and subroutines be as simple to use and as comprehensive in scope as is practical. Development of software that satisfies the ease of use criterion can be characterized as a packaging problem, the objective being to package mathematical theory and formulae into comprehensive, simple to use functions and subroutines. To help meet this criterion, many specialized routines are incorporated into the NSW library at a subordinate level, being referenced by driver functions or subroutines.

The testing of the functions and subroutines serves many purposes, including determination of the accuracy and efficiency of the software, checking for defects in the code, and searching for regions of numerical instability. In most cases, the testing must be highly



selective, being used to locate and examine weaknesses in the algorithm and code. After the testing is finished, an assessment is made of the utility and overall performance of the code. When the precision of a function or subroutine can be established, then this information is given with the description of the routine. All precision estimates are for the CDC 6000-7000 series 14-digit single and 28-digit double precision arithmetics. ***The estimates do not include inherent error.*** Thus, the accuracy of a code may be better than the inherent error of the mathematical function that it is computing.

The functions and subroutines in the NSWC library originate from a variety of sources. Hence, standards concerning in-line documentation and the style of code cannot be imposed. In general, all supportive routines not intended for general usage are not described in the manual. This makes it possible to modify or replace the code without bothering the programmer. This capability is extremely important. In the last decade, a vast amount of research has resulted in the development of new, more powerful algorithms for a variety of problems. Many of the results affect current codes, occasionally making some codes obsolete.

**Distribution of the NSWC library.** The NSWC functions and subroutines are available for general use. The library contains no proprietary code.

## MACHINE CONSTANTS

Assume that the integer arithmetic being used has base  $b$ , and that the integers are represented in the form

$$\pm(k_0 + k_1b + \dots + k_{n-1}b^{n-1})$$

where  $k_0, k_1, \dots, k_{n-1}$  are integers such that  $0 \leq k_i < b$  ( $i = 0, 1, \dots, n-1$ ). The value  $n$  is the number of base  $b$  digits  $k_i$ , and  $b^n - 1$  and  $-(b^n - 1)$  are the largest positive and negative integers that occur.

It is assumed that the single and double precision arithmetics being used have the same base, say  $\beta$ , and that the nonzero numbers can be represented in the form

$$\pm(\frac{k_1}{\beta} + \dots + \frac{k_m}{\beta^m})\beta^\ell$$

where  $k_1, \dots, k_m$  are integers such that  $0 \leq k_i < \beta$  ( $i = 1, \dots, m$ ),  $k_1 \geq 1$ , and  $\ell$  is an integer such that  $\ell_{\min} \leq \ell \leq \ell_{\max}$ . The value  $m$  is the number of base  $\beta$  digits  $k_i$ , and  $k_1 \geq 1$  is the requirement that the floating point numbers are normalized. The values  $\ell_{\min}$  and  $\ell_{\max}$  are the largest negative and positive exponents that arise where the floating point numbers maintain full  $m$  digit accuracy. Then  $x_{\min} = \beta^{\ell_{\min}-1}$  is the smallest positive number that is represented and  $x_{\max} = (1 - \beta^{-m})\beta^{\ell_{\max}}$  the largest.

Associated with  $m$  is the constant  $\epsilon = \beta^{-m+1}$ , called the *relative precision* of the floating point arithmetic being used. Theoretically,  $\epsilon$  is the smallest number for which  $1 + \epsilon$ , when stored in memory, is stored exactly and has a value greater than 1. Normally  $\epsilon$  will be the smallest number that satisfies these conditions. However, there do exist computer arithmetics (such as the CDC 6000-7000 series double precision arithmetics) for which this is not the case. In such arithmetics, some arithmetic results are able to be stored more accurately than others.

The functions SPMPAR, DPMPAR, and IPMPAR are available for obtaining the above constants. IPMPAR is the only subprogram in the NSWCLIB library that is machine dependent.

### SPMPAR( $i$ )

SPMPAR is a real valued function. It is assumed that  $i = 1, 2$ , or  $3$ . SPMPAR provides the following constants for the single precision arithmetic being used:

SPMPAR(1) =  $\epsilon$ , the relative precision

SPMPAR(2) =  $x_{\min}$ , the smallest positive number

SPMPAR(3) =  $x_{\max}$ , the largest positive number

**Programming.** SPMPAR was written by A. H. Morris. The function IPMPAR is used.

### DPMPAR(*i*)

DPMPAR is a double precision valued function. It is assumed that  $i = 1, 2$ , or  $3$ . DPMPAR provides the following constants for the double precision arithmetic being used:

- DPMPAR(1) =  $\epsilon$ , the relative precision
- DPMPAR(2) =  $x_{\min}$ , the smallest positive number
- DPMPAR(3) =  $x_{\max}$ , the largest positive number

DPMPAR must be declared in the calling program to be of type DOUBLE PRECISION.

**Programming.** DPMPAR was written by A. H. Morris. The function IPMPAR is used.

### IPMPAR(*i*)

IPMPAR is an integer valued function. It is assumed that  $i = 1, \dots, 10$ . IPMPAR provides the following constants:

#### *Integer arithmetic*

- IPMPAR(1) =  $b$ , the base of the arithmetic
- IPMPAR(2) =  $n$ , the number of base  $b$  digits
- IPMPAR(3) =  $b^n - 1$ , the largest integer

#### *Base for the floating arithmetics*

- IPMPAR(4) =  $\beta$

#### *Single precision arithmetic*

- IPMPAR(5) =  $m$ , the number of base  $\beta$  digits
- IPMPAR(6) =  $\ell_{\min}$ , the largest negative exponent
- IPMPAR(7) =  $\ell_{\max}$ , the largest positive exponent

#### *Double precision arithmetic*

- IPMPAR(8) =  $m$ , the number of base  $\beta$  digits
- IPMPAR(9) =  $\ell_{\min}$ , the largest negative exponent
- IPMPAR(10) =  $\ell_{\max}$ , the largest positive exponent

**Programming.** IPMPAR is an adaptation by A. H. Morris of the function IIMACH, designed by P. A. Fox, A. D. Hall, and N. L. Schryer (Bell Laboratories). The constants for the various computers are from Bell Laboratories, NSWC, and other sources.

## SORTING LISTS

Let  $A$  be an array containing  $n \geq 1$  elements  $a_1, \dots, a_n$ . Then the following subroutines are available for reordering the elements of  $A$ .

### CALL ISHELL( $A, n$ )

It is assumed that  $A$  is an integer array. When ISHELL is called, the elements of  $A$  are reordered so that  $a_i \leq a_{i+1}$  for  $i = 1, \dots, n-1$ .

**Algorithm.** The Shell sorting algorithm with increments  $(3^k - 1)/2$  is employed.

**Programmer.** A. H. Morris

**Reference.** Knuth, D. E., *The Art of Computer Programming. Vol. 3, Sorting and Searching*. Addison-Wesley, Reading, Mass., 1973, pp. 84-95.

### CALL SHELL( $A, n$ ) CALL AORD( $A, n$ )

It is assumed that  $A$  is a real array. If SHELL is called, then the elements of  $A$  are reordered so that  $a_i \leq a_{i+1}$  for  $i = 1, \dots, n-1$ . Otherwise, if AORD is called, then the elements of  $A$  are reordered so that  $|a_i| \leq |a_{i+1}|$  for  $i = 1, \dots, n-1$ .

**Programmer.** A. H. Morris

### CALL RISORT( $A, L, n$ )

It is assumed that  $A$  is a real array and  $L$  an integer array containing  $n$  elements. When RISORT is called, the elements of  $A$  are reordered so that  $a_i \leq a_{i+1}$  for  $i = 1, \dots, n-1$ . The same permutations are performed on  $L$  as on  $A$ , thereby reordering the elements of  $L$  so as to correspond with the new ordering of  $A$ .

**Remark.** RISORT is normally used for obtaining the indices of the reordered elements of  $A$ . If  $L$  initially contains the values  $1, \dots, n$  and  $a_{i_1}, \dots, a_{i_n}$  is the reordered sequence of elements of  $A$ , then  $L$  contains the indices  $i_1, \dots, i_n$  when RISORT terminates.

**Programmer.** A. H. Morris

### CALL SHELL2( $A, B, n$ )

It is assumed that  $A$  and  $B$  are real arrays containing  $n$  elements. When SHELL2 is called, the elements of  $A$  are reordered so that  $a_i \leq a_{i+1}$  for  $i = 1, \dots, n-1$ . The same permutations are performed on  $B$  as on  $A$ , thereby reordering the elements of  $B$  so as to correspond with the new ordering of  $A$ .

**Programmer.** A. H. Morris

**CALL DSORT( $A, n$ )**  
**CALL DAORD( $A, n$ )**

It is assumed that  $A$  is a double precision array. If DSORT is called, then the elements of  $A$  are reordered so that  $a_i \leq a_{i+1}$  for  $i = 1, \dots, n-1$ . Otherwise, if DAORD is called, then the elements of  $A$  are reordered so that  $|a_i| \leq |a_{i+1}|$  for  $i = 1, \dots, n-1$ .

**Programmer.** A. H. Morris

**CALL DISORT( $A, L, n$ )**

It is assumed that  $A$  is a double precision array and  $L$  an integer array containing  $n$  elements. When DISORT is called, the elements of  $A$  are reordered so that  $a_i \leq a_{i+1}$  for  $i = 1, \dots, n-1$ . The same permutations are performed on  $L$  as on  $A$ , thereby reordering the elements of  $L$  so as to correspond with the new ordering of  $A$ .

**Programmer.** A. H. Morris

**CALL DDSORT( $A, B, n$ )**

It is assumed that  $A$  and  $B$  are double precision arrays containing  $n$  elements. When DDSORT is called, the elements of  $A$  are reordered so that  $a_i \leq a_{i+1}$  for  $i = 1, \dots, n-1$ . The same permutations are performed on  $B$  as on  $A$ , thereby reordering the elements of  $B$  so as to correspond with the new ordering of  $A$ .

**Programmer.** A. H. Morris

## CUBE ROOT

The following functions are available for computing the real cube root of a real number.

**CBRT**( $x$ )  
**DCBRT**( $x$ )

CBRT is used if  $x$  is a single precision number, and DCBRT is used if  $x$  is a double precision number. CBRT is a single precision function and DCBRT a double precision function. The value of the function is  $\sqrt[3]{x}$ .

**Remark.** DCBRT must be declared in the calling program to be of type DOUBLE PRECISION.

**Programmer.** A. H. Morris

## FOUR QUADRANT ARCTANGENT

The function ARTNQ is similar to the ATAN2 function, the differences being that its value lies in the interval  $[0, 2\pi)$  and its value at the origin is 0. DARTNQ is the double precision counterpart of ARTNQ.

**ARTNQ**( $y, x$ )  
**DARTNQ**( $y, x$ )

ARTNQ is used if  $x$  and  $y$  are single precision values, and DARTNQ is used if  $x$  and  $y$  are double precision values. ARTNQ is a single precision function and DARTNQ is a double precision function.

If  $(x, y)$  is a point in the plane other than the origin  $(0, 0)$ , let  $L$  denote the straight line connecting the points  $(0, 0)$  and  $(x, y)$ . Then the function is assigned the value  $\theta$  where  $\theta$  is the angle between  $L$  and the positive  $x$ -axis measured in a counterclockwise direction. Here  $0 \leq \theta < 2\pi$ . Otherwise, if  $(x, y)$  is the origin  $(0, 0)$ , then the function is assigned the value 0.

**Remark.** DARTNQ must be declared in the calling program to be of type DOUBLE PRECISION.

**Programmer.** Richard Pasto

## LENGTH OF A TWO-DIMENSIONAL VECTOR

The following functions are available for computing the length of a real vector  $(x, y)$ .

**CPABS**( $x, y$ )  
**DCPABS**( $x, y$ )

CPABS is used if  $x$  and  $y$  are single precision values, and DCPABS is used if  $x$  and  $y$  are double precision values. CPABS is a single precision function and DCPABS is a double precision function. The value of the function is  $\sqrt{x^2 + y^2}$ .

**Remark.** DCPABS must be declared in the calling program to be of type DOUBLE PRECISION.

**Programmer.** A. H. Morris

## RECIPROCAL OF A COMPLEX NUMBER

The following subroutines are available for computing the reciprocal of a complex number.

**CALL CREC( $x, y, u, v$ )**  
**CALL DCREC( $x, y, u, v$ )**

CREC is used if  $x$  and  $y$  are single precision real numbers and  $u$  and  $v$  real variables, and DCREC is used if  $x$  and  $y$  are double precision numbers and  $u$  and  $v$  double precision variables. It is assumed that  $x$  and  $y$  are the real and imaginary parts of a nonzero complex number  $z$ . When CREC or DCREC is called,  $u$  and  $v$  are set to the real and imaginary parts of  $1/z$ , respectively.

**Programmer.** A. H. Morris

## SQUARE ROOT OF A DOUBLE PRECISION COMPLEX NUMBER

The following subroutine is available for computing the square root of a double precision complex number.

**CALL DCSQRT( $Z, W$ )**

$Z$  and  $W$  are double precision arrays of dimension 2. It is assumed that  $Z(1)$  and  $Z(2)$  are the real and imaginary parts of a complex number  $z$ . When DCSQRT is called, if  $z = 0$  then  $W(1)$  and  $W(2)$  are set to 0. Otherwise, if  $z \neq 0$  then the square root  $w = \sqrt{z}$  where  $-\pi/2 < \arg(w) \leq \pi/2$  is computed and stored in  $W$ .  $W(1)$  and  $W(2)$  contain the real and imaginary parts of  $w$ , respectively.

**Note.**  $Z$  and  $W$  may reference the same storage area.

**Programming.** DCSQRT calls the function DCPABS. DCSQRT was written by A. H. Morris.



## CONVERSION OF POLAR TO CARTESIAN COORDINATES

The following subroutine is available for converting polar coordinates  $(r, \theta)$  to cartesian coordinates  $(x, y)$ .

**CALL POCA** $(r, \theta, x, y)$

Let  $(r, \theta)$  be the polar coordinates of a point in the plane and let  $x, y$  be variables. When the routine is called,  $x$  and  $y$  are assigned the values  $x = r \cos \theta$  and  $y = r \sin \theta$ .

## CONVERSION OF CARTESIAN TO POLAR COORDINATES

The following subroutine is available for converting cartesian coordinates  $(x, y)$  to polar coordinates  $(r, \theta)$ .

**CALL CAPO** $(x, y, r, \theta)$

Let  $(x, y)$  be the cartesian coordinates of a point in the plane and let  $r, \theta$  be variables. If  $(x, y)$  is the origin then CAPO sets  $r = \theta = 0$ . Otherwise, if  $(x, y)$  is a point other than the origin, let  $L$  denote the straight line connecting the points  $(0, 0)$  and  $(x, y)$ . Then when CAPO is called,  $r$  is assigned the value  $\sqrt{x^2 + y^2}$  and  $\theta$  is defined to be the angle between  $L$  and the positive  $x$  axis. Here  $-\pi < \theta \leq \pi$ .

## ROTATION OF AXES

Let  $(x_1, y_1)$  be the (cartesian) coordinates for a point in the plane. The following subroutine computes the new coordinates  $(x_2, y_2)$  for the point after the  $x, y$  axes have been rotated by an angle  $\theta$ .

**CALL ROTA** $(x_1, y_1, \theta, x_2, y_2)$

The arguments  $x_2$  and  $y_2$  are variables. When ROTA is called,  $x_2$  and  $y_2$  are assigned the values:

$$\begin{aligned}x_2 &= x_1 \cos \theta + y_1 \sin \theta \\y_2 &= -x_1 \sin \theta + y_1 \cos \theta\end{aligned}$$

**Programmer.** A. H. Morris.

## PLANAR GIVENS ROTATIONS

If  $a$  and  $b$  are real numbers where  $a^2 + b^2 \neq 0$ , then there is an orthogonal matrix  $\begin{pmatrix} c & s \\ -s & c \end{pmatrix}$  such that  $\begin{pmatrix} c & s \\ -s & c \end{pmatrix} \begin{pmatrix} a \\ b \end{pmatrix} = \begin{pmatrix} r \\ 0 \end{pmatrix}$ . In this case  $r^2 = a^2 + b^2$ ,  $c = a/r$ , and  $s = b/r$ .

The matrix  $\begin{pmatrix} c & s \\ -s & c \end{pmatrix}$  represents what is called a *Givens rotation*. Given  $a$  and  $b$ , the matrix is uniquely defined up to the sign of  $r$ . For any real  $a$ , let  $\text{sgn}(a) = 1$  if  $a \geq 0$  and  $\text{sgn}(a) = -1$  if  $a < 0$ . If we define  $r = \sigma\sqrt{a^2 + b^2}$  where

$$\sigma = \begin{cases} \text{sgn}(a) & \text{if } |a| > |b| \\ \text{sgn}(b) & \text{if } |a| \leq |b| \end{cases}$$

then for  $r \neq 0$  we note that  $|c| > |s|$  implies  $c > 0$ , and that  $|c| \leq |s|$  implies  $s > 0$ . For convenience, when  $r = 0$  we set  $c = 1$  and  $s = 0$ .

The value  $\sigma$  is not needed for the construction of a Givens rotation matrix, but its use permits the representation of  $c$  and  $s$  by a single value  $z$ . For each  $c$  and  $s$ ,  $z$  is defined as follows:

$$z = \begin{cases} s & \text{if } |s| < c \text{ or } c = 0 \\ 1/c & \text{if } 0 < |c| \leq s \end{cases}$$

The mapping  $(c, s) \mapsto z$  is 1-1. If the user wishes to reconstruct  $c$  and  $s$  from  $z$ , then this can be done as follows:

If  $z = 1$  then set  $c = 0$  and  $s = 1$ .

If  $|z| < 1$  then set  $c = \sqrt{1 - z^2}$  and  $s = z$ .

If  $|z| > 1$  then set  $c = 1/z$  and  $s = \sqrt{1 - c^2}$ .

The subroutines SROTG and DROTG are available for computing  $c, s, r$ , and  $z$ . SROTG is used when  $a$  and  $b$  are single precision real numbers, and DROTG is used when  $a$  and  $b$  are double precision numbers.

**CALL SROTG(AR,BZ,C,S)**  
**CALL DROTG(AR,BZ,C,S)**

When SROTG is used then AR, BZ, C, and S are real variables. Otherwise, when DROTG is used then AR, BZ, C, and S are double precision variables. On input, AR =  $a$  and BZ =  $b$ . When the routine terminates AR =  $r$ , BZ =  $z$ , C =  $c$ , and S =  $s$ .

**Programming.** These routines are part of the BLAS package of basic linear algebra subroutines designed by C. L. Lawson, R. J. Hanson, D. R. Kincaid, and F. T. Krogh. SROTG and DROTG were coded by Charles Lawson (Jet Propulsion Laboratory).

### THREE DIMENSIONAL ROTATIONS

If  $A = (a_{ij})$  is a  $3 \times 3$  orthogonal matrix, then  $A$  can be represented in the form  $A = R_3 R_2 R_1 E$  where

$$R_1 = \begin{pmatrix} 1 & 0 & 0 \\ 0 & \cos \theta_1 & -\sin \theta_1 \\ 0 & \sin \theta_1 & \cos \theta_1 \end{pmatrix} \quad R_2 = \begin{pmatrix} \cos \theta_2 & 0 & -\sin \theta_2 \\ 0 & 1 & 0 \\ \sin \theta_2 & 0 & \cos \theta_2 \end{pmatrix}$$

$$R_3 = \begin{pmatrix} \cos \theta_3 & -\sin \theta_3 & 0 \\ \sin \theta_3 & \cos \theta_3 & 0 \\ 0 & 0 & 1 \end{pmatrix} \quad E = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & \pm 1 \end{pmatrix}.$$

$R_1$  represents a rotation around the  $x$ -axis,  $R_2$  a rotation around the  $y$ -axis, and  $R_3$  a rotation around the  $z$ -axis. Since  $A$  is orthogonal the determinant  $\det(A) = \pm 1$ . If  $\det(A) = 1$  then  $E$  is the identity matrix and  $A$  is the combined rotation  $R_3 R_2 R_1$ . If  $\det(A) = -1$  then  $A$  is composed of the rotation  $R_3 R_2 R_1$  and the reflection  $E = \text{diag}(1, 1, -1)$ . The following subroutine is available for finding the angles  $\theta_1, \theta_2, \theta_3$  where  $-\pi < \theta_1 \leq \pi, |\theta_2| \leq \pi/2$ , and  $-\pi < \theta_3 \leq \pi$ .

**CALL ROT3(A, THETA)**

THETA is an array of dimension 3 or larger. When ROT3 is called, the angles  $\theta_1, \theta_2, \theta_3$  are computed and stored in THETA.

**Algorithm.** If  $a_{11} = a_{21} = 0$  then let  $\theta_3 = 0$ . Otherwise, let  $\theta_3 = \text{ATAN2}(a_{21}, a_{11})$ . Then

$$R_3^t A = \begin{pmatrix} r_1 & a'_{12} & a'_{13} \\ 0 & a'_{22} & a'_{23} \\ a_{31} & a_{32} & a_{33} \end{pmatrix} = A'$$

where  $r_1 = \sqrt{a_{11}^2 + a_{21}^2}$ . Also, if  $\theta_2 = \text{ATAN2}(a_{31}, r_1)$  then

$$R_2^t A' = \begin{pmatrix} r_2 & a''_{12} & a''_{13} \\ 0 & a'_{22} & a'_{23} \\ 0 & a'_{32} & a'_{33} \end{pmatrix} = A''$$

where  $r_2 = \sqrt{r_1^2 + a_{31}^2}$ . Since  $r_2 \geq 0$ , by orthogonality it follows that  $r_2 = 1$  and  $a''_{12} = a''_{13} = 0$ . Finally, if  $\theta_1 = \text{ATAN2}(a'_{32}, a'_{22})$  then

$$R_1^t A'' = \begin{pmatrix} 1 & 0 & 0 \\ 0 & r_3 & a''_{23} \\ 0 & 0 & a''_{33} \end{pmatrix}$$

where  $r_3 = \sqrt{(a'_{22})^2 + (a'_{32})^2}$ . Since  $r_3 \geq 0$ , by orthogonality we obtain  $r_3 = 1$ ,  $a''_{23} = 0$ , and  $a''_{33} = \pm 1$ .

**Programmer.** A. H. Morris

## ROTATION OF A POINT ON THE UNIT SPHERE TO THE NORTH POLE

Given the point  $(x, y, z)$  where  $x^2 + y^2 + z^2 = 1$ . Then there exist orthogonal matrices

$$R_x = \begin{pmatrix} 1 & 0 & 0 \\ 0 & c_x & -s_x \\ 0 & s_x & c_x \end{pmatrix} \quad \text{and} \quad R_y = \begin{pmatrix} c_y & 0 & -s_y \\ 0 & 1 & 0 \\ s_y & 0 & c_y \end{pmatrix}$$

such that  $R_y R_x \begin{pmatrix} x \\ y \\ z \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix}$ .  $R_x$  represents a rotation about the  $x$ -axis and  $R_y$  a rotation about the  $y$ -axis. The following subroutine is available for obtaining the values  $c_x, s_x, c_y, s_y$ .

**CALL CONSTR**( $x, y, z, CX, SX, CY, SY$ )

$CX, SX, CY, SY$  are variables. When **CONSTR** is called, these variables are assigned the values  $c_x, s_x, c_y, s_y$ .

**Programmer.** Robert J. Renka (Oak Ridge National Laboratory)

## HYPERBOLIC SINE AND COSINE FUNCTIONS

The following subroutine is available for computing  $\sinh(x) - x$ ,  $\cosh(x) - 1$ , and  $\cosh(x) - 1 - x^2/2$  for real  $x$ .

**CALL SNHCSH( $S, C, x, \text{IND}$ )**

$S$  and  $C$  are variables, and  $\text{IND}$  is an input argument which specifies the functions to be computed.  $\text{IND}$  takes the values:

$\text{IND} = -1$  if only  $\sinh(x) - x$  is desired.

$\text{IND} = 0$  if  $\sinh(x) - x$  and  $\cosh(x) - 1$  are desired.

$\text{IND} = 1$  if only  $\cosh(x) - 1$  is desired.

$\text{IND} = 2$  if only  $\cosh(x) - 1 - x^2/2$  is desired.

$\text{IND} = 3$  if  $\sinh(x) - x$  and  $\cosh(x) - 1 - x^2/2$  are desired.

$S$  is assigned the value  $\sinh(x) - x$  if this function is requested. When  $\cosh(x) - 1$  or  $\cosh(x) - 1 - x^2/2$  is computed then the value is stored in  $C$ .

**Precision.** For all  $x$ ,  $\sinh(x) - x$  has a relative error less than  $2.3\text{E-}14$ ,  $\cosh(x) - 1$  has a relative error less than  $2.2\text{E-}14$ , and  $\cosh(x) - 1 - x^2/2$  has a relative error less than  $3.8\text{E-}14$ .

**Programming.** Written by A. K. Cline and R. J. Renka (University of Texas at Austin). Modified by A. H. Morris.

## EXPONENTIALS

The functions REXP and DREXP are available for computing  $e^x - 1$ . DREXP is a double precision function.

### REXP( $x$ )

$\text{REXP}(x) = e^x - 1$  for real  $x$ .

**Algorithm.** See pages 378-379 and appendix B of the reference.

**Precision.**  $\text{REXP}(x)$  is accurate to within 2 units of the 14<sup>th</sup> significant digit when  $\text{REXP}(x) \neq 0$ .

**Reference.** DiDonato, A. R., and Morris, A. H., "Computation of the Incomplete Gamma Function Ratios and Their Inverse," *ACM Trans. Math Software* 12 (1986), pp. 377-393.

**Programmer.** A. H. Morris

### DREXP( $x$ )

The argument  $x$  is a double precision number.  $\text{DREXP}(x)$  is the double precision value for  $e^x - 1$ .

**Remark.** DREXP must be declared in the calling program to be of type DOUBLE PRECISION.

**Precision.**  $\text{DREXP}(x)$  is accurate to within 1 unit of the 28<sup>th</sup> significant digit when  $\text{DREXP}(x) \neq 0$ .

**Programming.** DREXP was written by A. H. Morris. The function DPMPAR is used.

## LOGARITHMS

The functions ALNREL, RLOG, RLOG1, DLNREL, DRLOG, and DRLOG1 are available for computing  $\ln(1 + a)$ ,  $x - 1 - \ln(x)$ , and  $a - \ln(1 + a)$ . DLNREL, DRLOG, and DRLOG1 are double precision functions.

### ALNREL(*a*)

$$\text{ALNREL}(a) = \ln(1 + a) \text{ for } a > -1.$$

**Algorithm.** See page 378 and appendix A of the reference.

**Precision.** ALNREL(*a*) is accurate to within 2 units of the 14<sup>th</sup> significant digit when ALNREL(*a*)  $\neq$  0.

**Reference.** DiDonato, A. R., and Morris, A. H., "Computation of the Incomplete Gamma Function Ratios and Their Inverse," *ACM Trans. Math Software* 12 (1986), pp. 377-393.

**Programmer.** A. H. Morris

### RLOG(*x*)

$$\text{RLOG}(x) = x - 1 - \ln(x) \text{ for } x > 0.$$

**Algorithm.** See page 379 and appendix E of the reference.

**Precision.** RLOG(*x*) is accurate to within 3 units of the 14<sup>th</sup> significant digit when RLOG(*x*)  $\neq$  0.

**Reference.** DiDonato, A. R., and Morris, A. H., "Computation of the Incomplete Gamma Function Ratios and Their Inverse," *ACM Trans. Math Software* 12 (1986), pp. 377-393.

**Programmer.** A. H. Morris

### RLOG1(*a*)

$$\text{RLOG1}(a) = a - \ln(1 + a) \text{ for } a > -1.$$

**Algorithm.** See page 379 and appendix E of the reference.

**Precision.** RLOG1(*a*) is accurate to within 3 units of the 14<sup>th</sup> significant digit when RLOG1(*a*)  $\neq$  0.

**Reference.** DiDonato, A. R., and Morris, A. H., "Computation of the Incomplete Gamma Function Ratios and Their Inverse," *ACM Trans. Math Software* 12 (1986), pp. 377-393.

**Programmer.** A. H. Morris

### **DLNREL(*a*)**

The argument *a* is a double precision number where  $a > -1$ . DLNREL(*a*) is the double precision value for  $\ln(1 + a)$ .

**Remark.** DLNREL must be declared in the calling program to be of type DOUBLE PRECISION.

**Precision.** DLNREL(*a*) is accurate to within 1 unit of the 28<sup>th</sup> significant digit when DLNREL(*a*)  $\neq 0$ .

**Programming.** DLNREL was written by A. H. Morris. The function DPMPAR is used.

### **DRLOG(*x*)**

The argument *x* is a positive double precision number. DRLOG(*x*) is the double precision value for  $x - 1 - \ln(x)$ .

**Remark.** DRLOG must be declared in the calling program to be of type DOUBLE PRECISION.

**Precision.** DRLOG(*x*) is accurate to within 2 units of the 28<sup>th</sup> significant digit when DRLOG(*x*)  $\neq 0$ .

**Programming.** DRLOG was written by A. H. Morris. The function DPMPAR is used.

### **DRLOG1(*a*)**

The argument *a* is a double precision number  $a > -1$ . DRLOG1(*a*) is the double precision value for  $a - \ln(1 + a)$ .

**Remark.** DRLOG1 must be declared in the calling program to be of type DOUBLE PRECISION.

**Precision.** DRLOG1(*a*) is accurate to within 2 units of the 28<sup>th</sup> significant digit when DRLOG1(*a*)  $\neq 0$ .

**Programming.** DRLOG1 was written by A. H. Morris. The function DPMPAR is used.



## DETERMINING IF A POINT IS INSIDE OR OUTSIDE A POLYGON

Given a sequence of points  $\nu_i = (x_i, y_i)$  ( $i = 1, \dots, n$ ). Let  $\tau$  denote the polygonal line which begins at point  $\nu_1$ , traverses the points  $\nu_i$  in the order that they are indexed, and is the straight line segment connecting  $\nu_i$  to  $\nu_{i+1}$  for each  $i = 1, \dots, n-1$ . It is assumed that  $\nu_n = \nu_1$ , or that there is also a straight line segment from  $\nu_n$  to  $\nu_1$  when  $\nu_n \neq \nu_1$ . Consequently, the polygonal path  $\tau$  is a loop beginning and ending at  $\nu_1$ .

For any point  $\nu_0 = (x_0, y_0)$  not on the path  $\tau$ , let  $\eta(\tau, \nu_0)$  denote the *winding number* of the path  $\tau$  around the point  $\nu_0$ . Then  $\eta(\tau, \nu_0) = 0$  if  $\nu_0$  is outside the polygon whose boundary is  $\tau$ . If  $\nu_0$  is inside the path then  $\eta(\tau, \nu_0) = m$  where  $m$  is an integer. If  $m > 0$  then the path  $\tau$  loops  $m$  times in a counterclockwise direction around the point  $\nu_0$ . Otherwise, if  $m < 0$  then  $\tau$  loops  $|m|$  times in a clockwise direction around  $\nu_0$ .

Given an arbitrary point  $\nu_0 = (x_0, y_0)$ , the following subroutine is available for determining whether  $\nu_0$  is on the path, outside the path, or inside the path. If the point  $\nu_0$  is inside  $\tau$  then the winding number  $\eta(\tau, \nu_0) = m$  is also computed.

**CALL LOCPT**( $x_0, y_0, X, Y, n, \ell, m$ )

It is assumed that  $n \geq 1$ .  $X$  and  $Y$  are arrays containing  $x_1, \dots, x_n$  and  $y_1, \dots, y_n$ . The arguments  $\ell$  and  $m$  are variables. When LOCPT terminates  $\ell$  has one of the following values:

- $\ell = -1$  if  $(x_0, y_0)$  is outside the path
- $\ell = 0$  if  $(x_0, y_0)$  lies on the path
- $\ell = 1$  if  $(x_0, y_0)$  is inside the path

The variable  $m$  is assigned the value 0 if  $(x_0, y_0)$  is on or outside the path. If  $(x_0, y_0)$  is inside the path then  $m =$  the winding number of the path  $\tau$  around the point  $(x_0, y_0)$ .

**Remark.** There are no restrictions on the points  $(x_i, y_i)$ . Consequently, the path may intersect itself.

**Programming.** The function SPMPAR is used. LOCPT was written by A. H. Morris.

## THE CONVEX HULL FOR A FINITE PLANAR SET

If  $(x_1, y_1), \dots, (x_m, y_m)$  are  $m$  distinct points in the plane, then the following subroutine is available for finding the smallest convex polygon which with its interior contains the points.

**CALL HULL**( $X, Y, m, BX, BY, k, VX, VY, n$ )

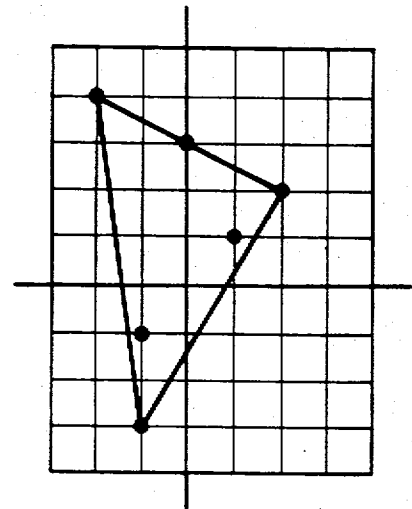
It is assumed that  $m \geq 2$ .  $X$  and  $Y$  are arrays containing the abscissas  $x_1, \dots, x_m$  and ordinates  $y_1, \dots, y_m$ , respectively. When HULL is called the points are reordered so that  $y_1 \leq \dots \leq y_m$ . Thus  $X$  and  $Y$  may be modified when the routine terminates.

$BX$  and  $BY$  are arrays of dimension  $m + 1$  or larger, and  $k$  is a variable. When HULL terminates,  $BX$  and  $BY$  contain the abscissas and ordinates of the points  $(x_i, y_i)$  which lie on the boundary of the desired convex polygon, and  $k$  = the number of points stored in  $BX$  and  $BY$ . If  $BX$  and  $BY$  contain the abscissas  $x'_1, \dots, x'_k$  and ordinates  $y'_1, \dots, y'_k$ , then the points  $(x'_i, y'_i)$  are indexed in the order they occur when traversing the boundary in a counterclockwise manner. Also  $(x'_k, y'_k) = (x'_1, y'_1)$ .

$VX$  and  $VY$  are arrays of dimension  $m + 1$  or larger, and  $n$  is a variable. When HULL terminates,  $VX$  and  $VY$  contain the abscissas and ordinates of the vertices of the desired convex polygon, and  $n$  = the number of points stored in  $VX$  and  $VY$ . If  $VX$  and  $VY$  contain the abscissas  $x''_1, \dots, x''_n$  and ordinates  $y''_1, \dots, y''_n$ , then the vertices  $(x''_i, y''_i)$  are indexed in the order they occur when traversing the boundary of the convex polygon in a counterclockwise manner. Also  $(x''_n, y''_n) = (x''_1, y''_1)$ .

**Example.** Assume that we are given the points  $(-1, -3), (1, 1), (0, 3), (2, 2), (-2, 4), (-1, -1)$ . When HULL is called  $X$  and  $Y$  are reordered and we obtain:

$X$  contains  $-1, -1, 1, 2, 0, -2$   
 $Y$  contains  $-3, -1, 1, 2, 3, 4$   
 $BX$  contains  $-1, 2, 0, -2, -1$   
 $BY$  contains  $-3, 2, 3, 4, -3$   
 $VX$  contains  $-1, 2, -2, -1$   
 $VY$  contains  $-3, 2, 4, -3$



**Programming.** HULL calls the subroutine RRSORT and function SPMPAR. HULL was written by A. H. Morris.

## AREAS OF PLANAR POLYGONS

Given a sequence of points  $v_i = (x_i, y_i)$  ( $i = 1, \dots, n+1$ ) where  $n \geq 1$  and  $v_{n+1} = v_1$ . Let  $\tau$  denote the polygon whose boundary  $\partial\tau$  is a polygonal line which begins at point  $v_1$ , traverses the points  $v_i$  in the order that they are indexed, and is the straight line segment connecting  $v_i$  to  $v_{i+1}$  for  $i = 1, \dots, n$ . Then the function PAREA is available for computing  $A(\tau) = \iint_{\tau} dx dy$ . If the boundary  $\partial\tau$  is a positively (negatively) oriented simple closed curve, then  $A(\tau)$  is positive (negative) and  $|A(\tau)|$  = the area of  $\tau$ . However,  $\partial\tau$  need not be simple. It may be self-intersecting or have overlapping line segments.

### PAREA( $X, Y, N$ )

$X$  and  $Y$  are arrays containing the abscissas  $x_1, \dots, x_N$  and ordinates  $y_1, \dots, y_N$ , respectively. The argument  $N$  may have the value  $n$  or  $n+1$ . Since  $v_{n+1} = v_1$ ,  $x_{n+1}$  and  $y_{n+1}$  are not required to appear in  $X$  and  $Y$ . PAREA( $X, Y, N$ ) is assigned the value  $A(\tau)$ .

**Programmer.** A. H. Morris

**Reference.** DiDonato, A. R. and Hageman, R. K., *Computation of the Integral of the Bivariate Normal Distribution over Arbitrary Polygons*, Report TR 80-166, Naval Surface Weapons Center, Dahlgren, Virginia, 1980.

## HAMILTONIAN CIRCUITS

Given a directed graph  $G$  containing  $n$  vertices, denoted by the integers  $1, \dots, n$ . Then any circuit of  $n$  arcs which traverses the  $n$  vertices, say in the order  $i_1, \dots, i_n$ , is called a **Hamiltonian circuit**. For convenience, it is assumed that for any two vertices  $i$  and  $j$ , no more than one arc exists which begins at  $i$  and ends at  $j$ . Then the following subroutine is available for finding the Hamiltonian circuits of  $G$ , if any exist.

**CALL HC(IND,  $m, n, P, A, NB, S, IWK, NUM$ )**

The argument  $m$  is the number of arcs in the graph  $G$ .  $P$  is an integer array of dimension  $n + 1$  and  $A$  an integer array of dimension  $m$ . The graph is stored in  $P$  and  $A$  as follows: For  $i = 1, \dots, n$  let

$$R_i = \{j: \text{there exists an arc which begins at } i \text{ and ends at } j\}.$$

Then the vertices in  $R_1, \dots, R_n$  are stored in  $A$ , where the data in  $R_i$  precedes the data in  $R_{i+1}$  for  $i = 1, \dots, n - 1$ . For each  $i$ , the vertices in  $R_i$  may be given in any order in  $A$ . The array  $P$  contains the data

$$P(1) = 0$$

$$P(i + 1) = \text{the total number of vertices in } R_1, \dots, R_i \text{ (} i = 1, \dots, n \text{)}.$$

Hence, if  $P(i) < P(i + 1)$  then the vertices in  $R_i$  are found in locations  $P(i) + 1, \dots, P(i + 1)$  of  $A$ . Otherwise, if  $P(i) = P(i + 1)$  then  $R_i = \phi$  (the empty set); i.e., there are no arcs in  $G$  which begin at  $i$  (and no Hamiltonian circuits exist). Also,  $P(n + 1) = m$  since there are  $m$  arcs in  $G$ .

**Example.** Consider the graph where

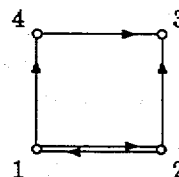
$$m = 5, n = 4$$

$$R_1 = \{2, 4\} \quad R_3 = \phi$$

$$R_2 = \{1, 3\} \quad R_4 = \{3\}.$$

Then  $A$  contains 4, 2, 1, 3, 3

$P$  contains 0, 2, 4, 4, 5.



When HC is called, a depth-first tree search employing backtracking is used.  $NB$  is a variable for controlling the backtracking. If  $NB \geq 0$  on input, then it is assumed that  $NB$  is the maximum number of backtracks that may be performed to find a Hamiltonian circuit. Otherwise, if  $NB < 0$  then it is assumed that no restriction is placed on the backtracking. When HC terminates,  $NB$  = the number of backtracks that were actually performed.

$S$  is an integer array of dimension  $n$ . When HC is called, if a Hamiltonian circuit is found which traverses the vertices, say in the order  $i_1, \dots, i_n$ , then the ordered vertices  $i_1, \dots, i_n$  are stored in  $S$ .

$IWK$  is an array of dimension  $NUM$  that is a work space for the routine. It is assumed that  $NUM \geq m + 8n + 20$ .

Only one Hamiltonian circuit will be obtained on a call to HC. However, this routine can be repeatedly called to obtain all the Hamiltonian circuits. IND is a variable which controls the operation of the routine on input, and reports the status of the results on output. It is assumed that  $IND = 0$  on the first call to HC. When the routine terminates, if no input errors are detected then IND has one of the following values:

IND = 1 A Hamiltonian circuit was found and the ordered vertices traversed by the circuit stored in  $S$ . To find another circuit, reset NB and recall the routine.

IND = 2 The maximum number of backtracks were performed. To continue, reset NB and recall the routine.

IND = 4 No more circuits exist. The array  $A$  has been restored (see the remark on the storage of  $A$  below) and the procedure is finished.

We note in passing that on an initial call to the routine, the setting  $IND = 4$  on output indicates that the graph contains no Hamiltonian circuits.

After a call to HC, if  $IND = 1$  or  $2$  on output then the search procedure can be continued by resetting NB and recalling the routine. *Do not modify IND when the tree search is to be continued.* In this case, the storage of  $A$  has been temporarily modified and IWK contains information needed for the search. If a new Hamiltonian circuit is found when HC is recalled, then the vertices traversed by the new circuit will now be stored in  $S$ .

After a call to HC, if  $IND = 1$  or  $2$  on output and it is desired that the search procedure be terminated, then reset  $IND = 3$  and recall the routine. In this case, the array  $A$  will be restored and  $IND = 4$  when HC terminates.

**Storage of  $A$ .** When  $A$  is restored, the order of the vertices of  $R_i$  in  $A$  may be modified for each  $i$ .

**Error Return.** If an input error is detected then IND is set to one of the following values:

IND = -1  $IND < 0$  or  $IND > 3$  on input.

IND = -2 IND was modified, being assigned a value  $\neq 3$  when HC was recalled. Reset IND to its previous output value, reset NB, and recall HC if another circuit is wanted.

IND = -3 The input setting  $IND = 3$  is not needed when the previous output value for IND was 4. In this case, nothing was done.

IND = -4  $NUM < m + 8n + 20$ .

IND = -5  $P(1) \neq 0$  or  $P(n+1) \neq m$ .

IND = -6  $P(i) > P(i+1)$  for some  $i$ .

**Remarks.**

- (1) It is assumed that  $G$  contains no loops.
- (2) Normally, few backtracks are needed when the number of vertices in each  $R_i$  is small, say 10 or less. Consequently, the setting  $NB = -1$  is generally appropriate in such cases.

**Programming.** HC employs the subroutines HC1, IPATH, FUPD, BUPD, IUPD, and RARC. The search procedure in these routines was written by Silvano Martello (University of Bologna, Italy). The user interface involving the variable IND was written by A. H. Morris.

**Reference.** Martello, S., "Algorithm 595. An Enumerative Algorithm for Finding Hamiltonian Circuits in a Directed Graph," *ACM Trans. Math Software* 9 (1983), pp. 131-138.

## ERROR FUNCTION

For any complex  $z$  the error function is defined by

$$\operatorname{erf}(z) = \frac{2}{\sqrt{\pi}} \int_0^z e^{-t^2} dt$$

and its complement by  $\operatorname{erfc}(z) = 1 - \operatorname{erf}(z)$ . The subroutines CERF, CERFC, DCERF, and DCERFC are available for computing  $\operatorname{erf}(z)$  and  $\operatorname{erfc}(z)$  when  $z$  is complex, and the functions ERF, ERFC, ERFC1, DERF, DERFC, and DERFC1 are available for computing  $\operatorname{erf}(z)$  and  $\operatorname{erfc}(z)$  when  $z$  is real. DCERF, DCERFC, DERF, DERFC, and DERFC1 are double precision routines.

### CALL CERF(MO, $z$ , $w$ )

MO is an integer,  $z$  a complex number, and  $w$  a complex variable. When CERF is called,  $w$  is assigned the value  $\operatorname{erf}(z)$  if MO = 0 and the value  $\operatorname{erfc}(z)$  if MO  $\neq$  0.

**Algorithm.** For  $z = x + iy$  where  $x \geq 0$ , if  $z$  satisfies  $|z| \leq 1$  or both of the inequalities  $1 < |z| < \sqrt{38}$  and  $x^2 - y^2 + .064x^2y^2 \leq 0$ , then the series

$$(1) \quad \operatorname{erf}(z) = \frac{2}{\sqrt{\pi}} \sum_{n \geq 0} \frac{(-1)^n z^{2n+1}}{n!(2n+1)}$$

is used. If  $1 < |z| < \sqrt{38}$  and  $x^2 - y^2 + .064x^2y^2 > 0$  then

$$(2) \quad \operatorname{erf}(z) = 1 - \frac{ze^{-z^2}}{\sqrt{\pi}} \sum_{n=1}^{18} \frac{r_n}{z^2 + \lambda_n}$$

is employed.  $\lambda_n$  and  $r_n$  are the poles and residues of the rational function approximation for the complex Fresnel integral  $E(z)$  given in the reference. The error function is related to  $E(z)$  by  $\operatorname{erf}(z) = 1 - i\sqrt{2}E(-z^2)$  for  $|\arg(z)| < \pi/2$ . If  $|z| \geq \sqrt{38}$  and  $x \geq .01$  then  $\operatorname{erf}(z)$  is computed by the asymptotic expansion  $\operatorname{erf}(z) = 1 - \psi(z)$  where

$$(3) \quad \psi(z) = \frac{e^{-z^2}}{\sqrt{\pi}} \left[ \frac{1}{z} + \sum_{n \geq 1} (-1)^n \frac{1 \cdot 3 \cdots (2n-1)}{2^n z^{2n+1}} \right].$$

Otherwise, if  $|z| \geq \sqrt{38}$  and  $0 \leq x < .01$  then  $\operatorname{erf}(z) = -\psi(z)$  is employed. When  $x < 0$  then the relation  $\operatorname{erf}(-z) = -\operatorname{erf}(z)$  is applied.

**Programming.** Written by Allen V. Hershey and A. H. Morris.

**Reference.** Hershey, A. V., *Approximation of Functions by Sets of Poles*, Report TR-2564, Naval Weapons Laboratory, Dahlgren, Virginia, 1971.

### CALL CERFC(MO, $z$ , $w$ )

MO is an integer,  $z$  a complex number, and  $w$  a complex variable. When CERFC is called,  $w$  is assigned the value  $\operatorname{erfc}(z)$  if  $\text{MO} = 0$  or  $\operatorname{Re}(z) < 0$ . Otherwise, if  $\text{MO} \neq 0$  and  $\operatorname{Re}(z) \geq 0$  then  $w$  is assigned the value  $e^{z^2}\operatorname{erfc}(z)$ .

**Precision.** For  $\text{MO} \neq 0$ ,  $\operatorname{Re}(w)$  and  $\operatorname{Im}(w)$  are accurate to within 1 unit of the  $12^{\text{th}}$  significant digit when  $|\operatorname{Re}(w)| \geq 10^{-280}$  and  $\operatorname{Im}(w) \neq 0$ .

**Programming.** CERFC employs the subroutine CREC and functions EXPARG and IPMPAR. CERFC was written by Allen V. Hershey and A. H. Morris.

**Reference.** Hershey, A. V., *Approximation of Functions by Sets of Poles*, Report TR-2564, Naval Weapons Laboratory, Dahlgren, Virginia, 1971.

### ERF( $x$ )

$\operatorname{ERF}(x) = \operatorname{erf}(x)$  for any real  $x$ .

**Precision.**  $\operatorname{ERF}(x)$  is accurate to within 2 units of the  $14^{\text{th}}$  significant digit for  $x \neq 0$ .

**Programmer.** A. H. Morris

**Reference.** Cody, W. J., "Rational Chebyshev Approximations for the Error Function," *Math Comp.* 23 (1969), pp. 631-637.

### ERFC( $x$ )

$\operatorname{ERFC}(x) = \operatorname{erfc}(x)$  for any real  $x$ .

**Precision.** If  $x \leq 3$  then  $\operatorname{ERFC}(x)$  is accurate to within 2 units of the  $14^{\text{th}}$  significant digit. Otherwise, if  $x > 3$  then  $\operatorname{ERFC}(x)$  is accurate to within 4 units of the  $14^{\text{th}}$  significant digit when  $\operatorname{ERFC}(x) \neq 0$ .

**Programmer.** A. H. Morris

**Reference.** Cody, W. J., "Rational Chebyshev Approximations for the Error Function," *Math Comp.* 23 (1969), pp. 631-637.

### ERFC1(IND, $x$ )

IND is an integer and  $x$  a real number.  $\operatorname{ERFC1}(\text{IND}, x) = \operatorname{erfc}(x)$  when  $\text{IND} = 0$ , and  $\operatorname{ERFC1}(\text{IND}, x) = e^{x^2}\operatorname{erfc}(x)$  when  $\text{IND} \neq 0$ .

**Precision.** If  $x \leq 3$  then  $\operatorname{ERFC1}(0, x)$  is accurate to within 2 units of the  $14^{\text{th}}$  significant digit. Otherwise, if  $x > 3$  then  $\operatorname{ERFC1}(0, x)$  is accurate to within 4 units of the  $14^{\text{th}}$  significant digit when  $\operatorname{ERFC1}(0, x) \neq 0$ . If  $\text{IND} \neq 0$  then  $\operatorname{ERFC1}(\text{IND}, x)$  is accurate to within 2 units of the  $14^{\text{th}}$  significant digit for  $x \geq -1$ .



**Programmer.** A. H. Morris

**Reference.** Cody. W. J., "Rational Chebyshev Approximations for the Error Function," *Math Comp.* 23 (1969), pp. 631-637.

**CALL DCERF(MO,Z,W)**

MO is an integer, and Z and W are double precision arrays of dimension 2. It is assumed that Z(1) and Z(2) are the real and imaginary parts of a complex number z. If MO = 0 then the double precision value for  $w = \operatorname{erf}(z)$  is computed, and if MO  $\neq$  0 then the double precision value for  $w = \operatorname{erfc}(z)$  is computed. W(1) and W(2) contain the real and imaginary parts of w, respectively.

**Algorithm.** For  $z = x + iy$  where  $x \geq 0$ , if  $|z| \leq 1$  then (1) is used, and if  $|z - 2| < 1$  and  $x \leq 2$  then the Taylor series

$$(4) \quad \operatorname{erfc}(z) = \operatorname{erfc}(a) + \frac{2}{\sqrt{\pi}} e^{-a^2} \sum_{n=1}^{\infty} (-1)^n H_{n-1}(a) (z-a)^n / n! \quad (a = 2)$$

is employed. Here  $H_n(a)$  are the Hermite polynomials. If  $1 < |z| < 2.5$  then (4) and the Pade approximation  $A_n(z^2)/B_n(z^2)$  for  $\sqrt{\pi}(2z)^{-1}e^{z^2}\operatorname{erf}(z)$  are used where

$$(5) \quad \begin{aligned} A_0(1) &= 1 & A_1(z) &= 1 + \frac{4}{15}z \\ B_0(1) &= 1 & B_1(z) &= 1 - \frac{2}{5}z \end{aligned}$$

and  $A_n$  and  $B_n$  satisfy

$$(6) \quad B_{n+1}(z) = \left[ 1 - \frac{2z}{(4n+1)(4n+5)} \right] B_n(z) + \frac{4n(4n+2)z^2}{(4n-1)(4n+1)^2(4n+3)} B_{n-1}(z)$$

(see pp. 191, 192, and 422 of the reference). Also, if  $2.5 \leq |z| < 12$  then  $A_n(z^2)/B_n(z^2)$  and the Pade approximation  $\overline{G}_n(z^2)/\overline{H}_n(z^2)$  for  $\sqrt{\pi}ze^{z^2}\operatorname{erfc}(z)$  are employed where

$$(7) \quad \begin{aligned} \overline{G}_0(z) &= 1 & \overline{G}_1(z) &= 2 + 2z \\ \overline{H}_0(z) &= 1 & \overline{H}_1(z) &= 3 + 2z \end{aligned}$$

and  $\overline{G}_n$  and  $\overline{H}_n$  satisfy

$$(8) \quad \overline{H}_{n+1}(z) = (2z + 4n + 3)\overline{H}_n(z) - 2n(2n+1)\overline{H}_{n-1}(z)$$

(see pp. 201 and 422 of the reference). Otherwise, if  $|z| \geq 12$  and  $x \geq .01$  then the asymptotic expansion  $\operatorname{erfc}(z) = \psi(z)$  is used where  $\psi(z)$  is given by (3). Also, if  $|z| \geq 12$  and  $0 \leq x < .01$  then  $\operatorname{erf}(z) = -\psi(z)$  is employed. When  $x < 0$  the relation  $\operatorname{erf}(-z) = -\operatorname{erf}(z)$  is applied.

**Programming.** DCERF calls the subroutines ERFCM2, CDIVID, and DCREC. DCERF

and ERFCM2 were written A. H. Morris. The function DPMPAR is also used.

**Reference.** Luke, Yudell L., *The Special Functions and Their Approximations, Vol 2*, Academic Press, New York, 1969.

### CALL DCERFC(MO,Z,W)

MO is an integer, and  $Z$  and  $W$  are double precision arrays of dimension 2. It is assumed that  $Z(1)$  and  $Z(2)$  are the real and imaginary parts of a complex number  $z$ . If  $MO = 0$  or  $\text{Re}(z) < 0$  then the double precision value for  $w = \text{erfc}(z)$  is computed. Otherwise, if  $MO \neq 0$  and  $\text{Re}(z) \geq 0$  then the double precision value for  $w = e^{z^2} \text{erfc}(z)$  is computed.  $W(1)$  and  $W(2)$  contain the real and imaginary parts of  $w$ , respectively.

**Precision.** For  $MO \neq 0$  and  $\text{Re}(z) \geq 0$ ,  $W(1)$  and  $W(2)$  are accurate to within 3 units of the  $26^{\text{th}}$  significant digit when  $W(2) \neq 0$ .

**Programming.** DCERFC employs the subroutines ERFCM2, DCIVID, and DCREC, and functions DXPAR, DPMPAR, and IPMPAR. DCERFC and ERFCM2 were written by A. H. Morris.

### DERF(x)

The argument  $x$  is a double precision real number. DERF( $x$ ) is the double precision value for  $\text{erf}(x)$ .

**Remark.** DERF must be declared in the calling program to be of type DOUBLE PRECISION.

**Algorithm.** If  $|x| \leq 1$  then the power series corresponding to the Chebyshev expansion given in the SLATEC library by Wayne Fullerton (Los Alamos) is used. The power series was obtained by A. H. Morris. If  $|x| > 1$  then Chebyshev expansions derived by J. L. Schonfelder (University of Birmingham, England) are used.

**Precision.** DERF( $x$ ) is accurate to within 2 units of the  $28^{\text{th}}$  significant digit for  $x \neq 0$ .

**Programming.** DERF calls the functions DCSEVL and DPMPAR. DERF was written by A. H. Morris.

**Reference.** Schonfelder, J. L., "Chebyshev Expansions for the Error and Related Functions," *Math Comp.* 32 (1978), pp. 1232-1240.

### DERFC(x)

The argument  $x$  is a double precision real number. DERFC( $x$ ) is the double precision value for  $\text{erfc}(x)$ .

**Remark.** DERFC must be declared in the calling program to be of type DOUBLE PRECISION.

**Algorithm.** If  $|x| \leq 1$  then the power series corresponding to the Chebyshev expansion given in the SLATEC library by Wayne Fullerton (Los Alamos) is used. The power series was obtained by A. H. Morris. If  $|x| > 1$  then Chebyshev expansions derived by J. L. Schonfelder (University of Birmingham, England) and the standard asymptotic expansion for  $\operatorname{erfc}(x)$  are used.

**Precision.**  $\operatorname{DERFC}(x)$  is accurate to within 2 units of the  $28^{\text{th}}$  significant digit for  $x \leq 2$ .

**Programming.**  $\operatorname{DERFC}$  calls the functions  $\operatorname{DCSELV}$  and  $\operatorname{DPMPAR}$ .  $\operatorname{DERFC}$  was written by A. H. Morris.

**Reference.** Schonfelder, J. L., "Chebyshev Expansions for the Error and Related Functions," *Math Comp.* 32 (1978), pp. 1232-1240.

### $\operatorname{DERFC1}(\operatorname{IND}, x)$

$\operatorname{IND}$  is an integer and  $x$  a double precision real number.  $\operatorname{DERFC1}(\operatorname{IND}, x)$  is the double precision value for  $\operatorname{erfc}(x)$  when  $\operatorname{IND} = 0$ , and  $\operatorname{DERFC1}(\operatorname{IND}, x)$  is the double precision value for  $e^{x^2}\operatorname{erfc}(x)$  when  $\operatorname{IND} \neq 0$ .

**Remark.**  $\operatorname{DERFC1}$  must be declared in the calling program to be of type **DOUBLE PRECISION**.

**Precision.**  $\operatorname{DERFC1}(\operatorname{IND}, x)$  is accurate to within 2 units of the  $28^{\text{th}}$  significant digit when  $\operatorname{IND} = 0$  and  $x \leq 2$ , and when  $\operatorname{IND} \neq 0$  and  $x \geq -1$ .

**Programming.**  $\operatorname{DERFC1}$  calls the functions  $\operatorname{DCSEVL}$  and  $\operatorname{DPMPAR}$ .  $\operatorname{DERFC1}$  was written by A. H. Morris.

## INVERSE ERROR FUNCTION

For any  $0 \leq x < 1$ , the following function is available for obtaining the value  $w \geq 0$  for which  $\text{erf}(w) = x$ .

**ERFINV**( $x, y$ )

It is assumed that  $0 \leq x \leq 1$  and  $y = 1 - x$ . If  $y \neq 0$  then  $\text{ERFINV}(x, y) = w$  where  $w \geq 0$  and  $\text{erf}(w) = x$ . Otherwise, if  $y = 0$  then  $\text{ERFINV}(x, y) =$  the largest positive number in the floating arithmetic being used.

**Error Return.**  $\text{ERFINV}(x, y) < 0$  if  $x < 0$ ,  $y < 0$ , or  $x + y \neq 1$ .

**Precision.** For  $x \neq 0$  and  $y \neq 0$ ,  $\text{ERFINV}(x, y)$  is accurate to within 3 units of the 14<sup>th</sup> significant digit.

**Programming.**  $\text{ERFINV}$  was written by Armido R. DiDonato and modified by A. H. Morris. The function  $\text{SPMPAR}$  is used.

**Reference.** Blair, J. M., Edwards, C. A., and Johnson, J. H., "Rational Chebyshev Approximations for the Inverse of the Error Function," *Math Comp.* 30 (1976), pp. 827-830.

## NORMAL PROBABILITY DISTRIBUTION FUNCTION

For any real  $x$ , the normal probability distribution function  $P(x)$  (of mean 0 and variance 1) is defined by

$$P(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-t^2/2} dt$$

and its complement by  $Q(x) = 1 - P(x)$ . The following function is available for computing  $P(x)$  and  $Q(x)$ .

**PNDF**( $x$ ,IND)

IND is an integer and  $x$  a real number. If IND = 0 then

$$\text{PNDF}(x,0) = \begin{cases} P(x) & \text{if } x \geq -8 \\ \frac{P'(x)}{P(x)} & \text{if } x < -8 \end{cases}$$

where  $P'(x)$  is the derivative of  $P(x)$ . Otherwise, if IND  $\neq$  0 then

$$\text{PNDF}(x,\text{IND}) = \begin{cases} Q(x) & \text{if } x \leq 8 \\ \frac{Q'(x)}{Q(x)} & \text{if } x > 8 \end{cases}$$

where  $Q'(x)$  is the derivative of  $Q(x)$ .

**Algorithm.** The identities  $P(x) = \frac{1}{2} \text{erfc}(-x/\sqrt{2})$  and  $Q(x) = \frac{1}{2} \text{erfc}(x/\sqrt{2})$  are used.

**Programming.** PNDF calls the function ERFC1. PNDF was written by A.H. Morris.

## INVERSE NORMAL PROBABILITY DISTRIBUTION FUNCTION

For any real  $w$ , the normal probability distribution function  $P(w)$  (of mean 0 and variance 1) is defined by

$$P(w) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^w e^{-t^2/2} dt$$

and its complement by  $Q(w) = 1 - P(w)$ . For any  $0 < p < 1$  and  $q = 1 - p$ , the following function is available for obtaining the value  $w$  for which  $P(w) = p$  and  $Q(w) = q$ .

**PNINV( $p, q, z, \text{IERR}$ )**

It is assumed that  $0 \leq p \leq 1, q = 1 - p$ , and  $z = p - q$ . IERR is a variable. When PNINV is used, if no input errors are detected then IERR is set to 0. If  $p \neq 0$  and  $q \neq 0$  then  $\text{PNINV}(p, q, z, \text{IERR}) = w$  where  $P(w) = p$  and  $Q(w) = q$ . Otherwise,

$$\text{PNINV}(p, q, z, \text{IERR}) = \begin{cases} -x_{\max} & \text{if } p = 0 \\ x_{\max} & \text{if } q = 0 \end{cases}$$

where  $x_{\max}$  is the largest positive number in the floating arithmetic being used.

**Error Return.** If an input error is detected, then IERR is set to one of the following values:

IERR = 1 Either  $p < 0, q < 0$ , or  $p + q \neq 1$ .

IERR = 2  $z \neq p - q$

In these cases, PNINV is assigned the value 0.

**Algorithm.** For  $y \geq 0$ , let  $y = \text{erf}^{-1}(x)$  when  $x = \text{erf}(y)$ . If  $P(w) = p$  then the identities

$$w = \begin{cases} -\sqrt{2} \text{erf}^{-1}(1 - 2p) & \text{if } 0 < p \leq 1/2 \\ \sqrt{2} \text{erf}^{-1}(2p - 1) & \text{if } 1/2 < p \leq 1 \end{cases}$$

are applied.

**Precision.** For  $p \neq 0, q \neq 0$ , and  $z \neq 0$  PNINV is accurate to within 3 units of the 14<sup>th</sup> significant digit.

**Programming.** PNINV employs the functions ERFINV and SPMPAR. PNINV was written by Armido R. DiDonato.

## DAWSON'S INTEGRAL

For any real  $x$ , Dawson's integral is defined by

$$F(x) = e^{-x^2} \int_0^x e^{t^2} dt.$$

The following function is available for computing  $F(x)$ .

**DAW( $x$ )**

DAW( $x$ ) =  $F(x)$  for any real  $x$ .

**Precision.** DAW( $x$ ) is accurate to within 2 units of the 14<sup>th</sup> significant digit for  $x \neq 0$ .

**Programming.** DAW belongs to the FUNPACK package of subroutines developed at Argonne National Laboratory. The function was modified by A. H. Morris.

**Reference.** Cody, W. J., Paciorek, K. A., and Tacher, H. C., "Chebyshev Approximations for Dawson's Integral," *Math Comp.* 24 (1970), pp. 171-178.

## COMPLEX FRESNEL INTEGRAL

For any complex  $z$  not on the positive real axis the complex Fresnel integral  $E(z)$  can be defined by

$$E(z) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^z \frac{e^t}{\sqrt{t}} dt.$$

Here it is assumed that  $0 < \arg(z) < 2\pi$  and  $\arg(\sqrt{z}) = 1/2 \arg(z)$ .  $E(z)$  can be extended to the positive real axis by letting  $0 \leq \arg(z) < 2\pi$ . Then  $\operatorname{erf}(z) = 1 - i\sqrt{2} E(-z^2)$  for  $-\pi/2 \leq \arg(z) < \pi/2$  where  $\operatorname{erf}(z)$  is the error function and  $\arg(-z^2) = \pi + 2 \arg(z)$ . The following subroutine is available for computing  $E(z)$ .

**CALL CFRNLI(MO,  $z$ ,  $w$ )**

MO is an integer,  $z$  a complex number, and  $w$  a complex variable. When CFRNLI is called,  $w$  is assigned the value  $E(z)$  if  $\text{MO} = 0$  and the value  $e^{-z} E(z)$  if  $\text{MO} \neq 0$ .

**Algorithm.** If  $z = x + iy$  satisfies  $|z| \leq 1$  or both of the inequalities  $1 < |z| < 38$  and  $-x + .016y^2 \leq 0$ , then the series

$$E(z) = -\frac{i}{\sqrt{2}} + \sqrt{2z/\pi} \sum_{n \geq 0} \frac{z^n}{n!(2n+1)}$$

is used. If  $1 < |z| < 38$  and  $-x + .016y^2 > 0$  then

$$E(z) = \sqrt{\frac{z}{2\pi}} e^z \sum_{n=1}^{18} \frac{r_n}{z - \delta_n}$$

is employed. If  $|z| \geq 38$  and  $\operatorname{Im} \sqrt{2z/\pi} \geq .008$  then  $E(z)$  is computed by the asymptotic expansion  $E(z) = \psi(z)$  where

$$\psi(z) = \frac{e^z}{\sqrt{2\pi z}} \left[ 1 + \sum_{n \geq 1} \frac{1 \cdot 3 \cdots (2n-1)}{(2z)^n} \right].$$

Otherwise, if  $|z| \geq 38$  and  $\operatorname{Im} \sqrt{2z/\pi} < .008$  then  $E(z) = -i/\sqrt{2} + \psi(z)$  is employed.

**Precision.** For  $\text{MO} \neq 0$ ,  $\operatorname{Re}(w)$  and  $\operatorname{Im}(w)$  are accurate to within 1 unit of the 12<sup>th</sup> significant digit when  $\operatorname{Re}(w)$  and  $\operatorname{Im}(w)$  are nonzero.

**Programming.** CFRNLI employs the functions CPABS, EXPARG, and IPMPAR. CFRNLI was written by Allen V. Hershey and A. H. Morris.

**Reference.** Hershey, A. V., *Approximation of Functions by Sets of Poles*, Report TR-2564, Naval Weapons Laboratory, Dahlgren, Virginia, 1971.



## REAL FRESNEL INTEGRALS

For any complex  $z$  the Fresnel integrals  $C(z)$  and  $S(z)$  can be defined by

$$C(z) = \int_0^z \cos\left(\frac{\pi}{2}t^2\right) dt$$

$$S(z) = \int_0^z \sin\left(\frac{\pi}{2}t^2\right) dt.$$

The following subroutine is available for computing  $C(z)$  and  $S(z)$  when  $z$  is real.

**CALL FRNL**( $x, C, S$ )

The argument  $x$  may be any real number.  $C$  and  $S$  are variables. When FRNL is called  $C$  is assigned the value  $C(x)$  and  $S$  is assigned the value  $S(x)$ .

**Algorithm.** If  $0 \leq x < 1.65$  then  $x^{-1}C(x)$  and  $x^{-3}S(x)$  are computed by minimax polynomial approximations. Otherwise, if  $x \geq 1.65$  then the relations

$$C(x) = \frac{1}{2} + f(x) \sin \frac{\pi}{2}x^2 - g(x) \cos \frac{\pi}{2}x^2$$

$$S(x) = \frac{1}{2} - f(x) \cos \frac{\pi}{2}x^2 - g(x) \sin \frac{\pi}{2}x^2$$

are invoked. For  $1.65 \leq x < 6$ ,  $xf(x)$  and  $x^3g(x)$  are computed by rational minimax approximations. Otherwise, for  $x \geq 6$  the auxiliary functions  $f(x)$  and  $g(x)$  are computed by the asymptotic expansions:

$$f(x) = \frac{1}{\pi x} \left[ 1 + \sum_{i=1}^m (-1)^i \frac{1 \cdot 3 \cdots (4i-1)}{(\pi x^2)^{2i}} \right]$$

$$g(x) = \frac{1}{\pi x} \sum_{i=0}^m (-1)^i \frac{1 \cdot 3 \cdots (4i+1)}{(\pi x^2)^{2i+1}}$$

Here  $m = 5$ . If  $x < 0$  then the relations  $C(-x) = -C(x)$  and  $S(-x) = -S(x)$  are applied.

**Precision.** If  $|x| < 1.65$  then FRNL is accurate to within 3 units of the 14<sup>th</sup> significant digit. Otherwise, if  $|x| \geq 1.65$  then FRNL is accurate to within 1 unit of the 14<sup>th</sup> significant digit.

**Programming.** FRNL calls the functions SPMPAR and IPMPAR. FRNL was written by A. H. Morris.

## EXPONENTIAL INTEGRAL FUNCTION

For any complex  $z \neq 0$  not on the positive real axis the exponential integral function  $Ei(z)$  is defined by

$$Ei(z) = \int_{-\infty}^z \frac{e^t}{t} dt.$$

$Ei(z)$  is an analytic function. If  $z$  is replaced by  $-z$  and  $t$  by  $-t$  we obtain the related function

$$E_1(z) = -Ei(-z) = \int_z^{\infty} \frac{e^{-t}}{t} dt$$

which is defined for all  $z \neq 0$  not on the negative real axis. It can be verified that

$$Ei(z) = \nu + \ln(-z) + \sum_{n=1}^{\infty} \frac{z^n}{n \cdot n!}$$

everywhere in the plane cut along the positive real axis where  $\nu$  is the Euler constant. Thus, the values of  $Ei(x)$  on the upper and lower edges of the cut are

$$Ei(x \pm i0) = ei(x) \mp \pi i$$

where  $ei(x)$  is the real function defined by

$$ei(x) = \nu + \ln x + \sum_{n=1}^{\infty} \frac{x^n}{n \cdot n!}$$

for  $x > 0$ . The function  $ei(x)$ , also known as the exponential integral function, has a zero at the point  $x_0 = .37250\ 74107\ 81367$ .  $Ei(z)$  may be computed by the subroutine CEXPLI when  $z$  is complex, and  $Ei(z)$  and  $ei(z)$  may be computed by the subroutine EXPLI and functions DEI and DEI1 when  $z$  is real. DEI and DEI1 are double precision functions.

### CALL CEXPLI(MO, z, w)

MO is an integer,  $z \neq 0$  a complex number, and  $w$  a complex variable. When CEXPLI is called,  $w$  is assigned the value  $Ei(z)$  if  $MO = 0$  and the value  $e^{-z} Ei(z)$  if  $MO \neq 0$ .

**Remark.** If  $z$  is a positive real number and  $MO = 0$  then  $w = ei(z) + \pi i$ .

**Precision.** If  $MO = 0$  then  $Re(w)$  and  $Im(w)$  are accurate to within 2 units of the 12<sup>th</sup> significant digit when  $z$  is not near a zero of  $Re(Ei(z))$  or  $Im(Ei(z))$ .

**Programming.** CEXPLI employs the functions CPABS and SPMPAR. CEXPLI was initially written by Allen V. Hershey, and later rewritten by A. H. Morris.

**Reference.** Hershey, A. V., *Approximations of Functions by Sets of Poles*, Report TR-2564, Naval Weapons Laboratory, Dahlgren, Virginia, 1971.

## CALL EXPLI(MO, $x$ , $w$ , IERR)

MO may have the values 1, 2, or 3. The argument  $x$  is a nonzero real number and  $w$  a real variable. When EXPLI is called, if MO = 1 then  $w$  is assigned the value  $Ei(x)$  for  $x < 0$  and the value  $ei(x)$  for  $x > 0$ . If MO = 2 then it is assumed that  $x > 0$ . In this case  $w$  is assigned the value  $E_1(x)$ . Otherwise, if MO = 3 then  $w$  is assigned the value  $e^{-x}Ei(x)$  for  $x < 0$  and the value  $e^{-x}ei(x)$  for  $x > 0$ .

**Error return.** IERR is a variable that reports the status of the results. If the requested value  $w$  is obtained then IERR is set to 0. Otherwise, IERR is assigned one of the following values:

IERR = 1 Underflow occurs. In this case  $w = 0$ .

IERR = 2 Overflow occurs.

IERR = 3 (Input error)  $x = 0$ .

IERR = 4 (Input error) MO = 2 and  $x < 0$ .

The variable  $w$  is not defined when IERR  $\geq 2$ .

**Algorithm.** If MO  $\neq 2$  and  $4 \leq x \leq 8$ , then the Chebyshev expansion in the SLATEC library obtained by Wayne Fullerton (Los Alamos) is used. The remaining approximations employed are from the references.

**Precision.** If MO  $\neq 2$  and  $x > 0$ , then  $w$  is accurate to within 4 units of the 14<sup>th</sup> significant digit when  $w \neq 0$ . Otherwise,  $w$  is accurate to within 3 units of the 14<sup>th</sup> significant digit when  $w = 0$ .

**Programming.** EXPLI employs the functions ALNREL, CSEVL, EXPARG, and IPMPAR. EXPLI was written at Argonne National Laboratory for the FUNPACK package of special function subroutines. EXPLI was modified by A. H. Morris.

## References

- (1) Cody, W. J. and Thacher, H. C., "Rational Chebyshev Approximations for the Exponential Integral  $E_1(x)$ ," *Math Comp.* 22 (1968), pp. 641-649.
- (2) \_\_\_\_\_, "Chebyshev Approximations for the Exponential Integral  $Ei(x)$ ," *Math Comp.* 23 (1969), pp. 289-303.

## DEI( $x$ )

The argument  $x \neq 0$  is a double precision number. DEI( $x$ ) is the double precision value for  $Ei(x)$  when  $x < 0$ , and the double precision value for  $ei(x)$  when  $x > 0$ .

**Remark.** DEI must be declared in the calling program to be of type DOUBLE PRECISION.

**Algorithm.** If  $.35 \leq x \leq .4$  then the Taylor series expansion of  $ei(x)$  around  $x_0$  is used. This expansion was obtained by A. H. Morris. If  $|x| \geq 90$  then the standard asymptotic expansion for  $Ei$  and  $ei$  is applied. Otherwise, the Chebyshev expansions in the SLATEC library obtained by Wayne Fullerton (Los Alamos) are used.

**Programming.** DEI employs the functions DE1E and DEI0. These functions were written by A. H. Morris. The functions DCSEVL and DPMPAR are also used.

### **DEI1( $x$ )**

The argument  $x \neq 0$  is a double precision number. DEI1( $x$ ) is the double precision value for  $e^{-x} \text{Ei}(x)$  when  $x < 0$ , and the double precision value for  $e^{-x} \text{ei}(x)$  when  $x > 0$ .

**Remark.** DEI1 must be declared in the calling program to be of type DOUBLE PRECISION.

**Algorithm.** If  $.35 \leq x \leq .4$  then the Taylor series expansion of  $\text{ei}(x)$  around  $x_0$  is used. This expansion was obtained by A. H. Morris. If  $|x| \geq 90$  then the standard asymptotic expansion for Ei and ei is applied. Otherwise, the Chebyshev expansions in the SLATEC library obtained by Wayne Fullerton (Los Alamos) are used.

**Precision.** DEI1( $x$ ) is accurate to within 4 units of the 28<sup>th</sup> significant digit when DEI1( $x$ )  $\neq 0$ .

**Programming.** DEI1 employs the functions DE1E and DEI0. These functions were written by A. H. Morris. The functions DCSEVL and DPMPAR are also used.

## SINE AND COSINE INTEGRAL FUNCTIONS

For any complex  $z$  the sine integral and cosine integral functions  $\text{Si}(z)$  and  $\text{Cin}(z)$  are defined by

$$\begin{aligned}\text{Si}(z) &= \int_0^z \frac{\sin t}{t} dt \\ \text{Cin}(z) &= \int_0^z \frac{1 - \cos t}{t} dt.\end{aligned}$$

These are entire functions. The following functions are available for computing  $\text{Si}(z)$  and  $\text{Cin}(z)$  when  $z$  is real.

### **SI( $x$ )**

$\text{SI}(x) = \text{Si}(x)$  for all real  $x$ .

**Precision.** SI is accurate to within 2 units of the 14<sup>th</sup> significant digit.

**Programming.** SI calls the function SPMPAR. SI was written by Donald E. Amos and Sharon L. Daniel (Sandia Laboratories), and modified by A. H. Morris.

### **CIN( $x$ )**

$\text{CIN}(x) = \text{Cin}(x)$  for all real  $x$ .

**Precision.** CIN is accurate to within 2 units of the 14<sup>th</sup> significant digit.

**Programming.** CIN calls the function SPMPAR. CIN was written by Donald E. Amos and Sharon L. Daniel (Sandia Laboratories), and modified by A. H. Morris.

## DILOGARITHM FUNCTION

For any complex  $z$  where  $|\arg(1+z)| < \pi$ , the dilogarithm function  $L(z)$  may be defined by

$$L(z) = \int_0^z \frac{\ln(1+t)}{t} dt.$$

$L(z)$  is real-valued for any real  $z \geq -1$ , and  $-1$  is a branch point.  $L(z)$  can be extended to the negative real axis from  $-\infty$  to  $-1$  by letting  $-\pi < \arg(1+z) \leq \pi$ . Then for any real  $x < -1$

$$L(x) = -\pi^2/6 + \int_1^{-x} \frac{\ln(t-1)}{t} dt + i\pi \ln(-x) \quad (i = \sqrt{-1}).$$

The function CLI is available for computing  $L(z)$  when  $z$  is complex, and the function ALI is available for computing the real part of  $L(z)$  when  $z$  is real.

### CLI( $z$ )

CLI is a complex-valued function where  $\text{CLI}(z) = L(z)$  for all complex  $z$ . CLI must be declared in the calling program to be of type COMPLEX.

**Algorithm.** For  $|z| \leq 1/2$  the Maclaurin series

$$(1) \quad L(z) = - \sum_{n \geq 1} \frac{(-z)^n}{n^2}$$

is used. If  $|z| \geq 3$  then

$$(2) \quad L(z) = \pi^2/6 - L(1/z) + 1/2 \ln^2 z$$

is applied, and if  $0 < |z+1| \leq 1/2$  then

$$(3) \quad L(z) = -\pi^2/6 - L(-1-z) + \ln(-z) \ln(1+z)$$

is applied. Otherwise,

$$(4) \quad \begin{aligned} L(z) &= - \int_0^w \frac{\lambda}{e^\lambda - 1} d\lambda \quad (\text{Debye Function}) \\ &= -w + w^2/4 - \sum_{n \geq 1} \frac{B_{2n} w^{2n+1}}{(2n+1)!} \quad |w| < 2\pi \end{aligned}$$

is used where  $w = -\ln(1+z)$  and  $B_{2n}$  are the Bernoulli numbers  $B_2 = 1/6$ ,  $B_4 = -1/30$ , .... In (3) we note that  $\ln(-z) \ln(1+z) \rightarrow 0$  when  $z \rightarrow -1$ .

**Programming.** CLI is a modification by A. H. Morris of the subroutine CLGMCI, written by Allen V. Hershey.

**Reference.** Hershey, A. V., *Approximation of Functions by Sets of Poles*, Report TR-2564, Naval Weapons Laboratory, Dahlgren, Virginia, 1971.

**ALI( $x$ )**

$\text{ALI}(x) = \text{Re}[L(x)]$  for all real  $x$ .

**Algorithm.** Rational minimax approximations are used when  $-1/2 \leq x \leq 1$ . If  $x > 1$  then (2) is applied, and if  $-2 \leq x < -1$  or  $-1 < x < -1/2$  then (3) is applied. Otherwise, if  $x < -2$  then

$$\text{Re}[L(x)] = -\pi^2/3 - L(1/x) + 1/2 \ln^2(-x)$$

(which follows from (2)) is used except when  $-26.63 \leq x \leq -6.97$ . Since  $\text{Re}[L(x_0)] = 0$  for  $x_0 = -12.59517 \dots$ , the Taylor series of  $\text{Re}[L(x)]$  around  $x_0$  is used when  $-14 < x \leq -11.1$ . Otherwise, if  $-11.1 < x \leq -6.97$  or  $-26.63 \leq x \leq -14$  then rational minimax approximations are employed.

**Precision.**  $\text{ALI}(x)$  is accurate to within 2 units of the 14<sup>th</sup> significant digit when  $x > 0$ . Otherwise, if  $x < 0$  then  $\text{ALI}(x)$  is accurate to within 4 units of the 14<sup>th</sup> significant digit when  $\text{ALI}(x) \neq 0$ .

**Programmer.** A. H. Morris

**Reference.** Morris, Robert, "The Dilogarithm Function of a Real Argument," *Math. Comp.* **33** (1979), pp. 778-787.

## GAMMA FUNCTION

For any complex  $z \neq 0, -1, -2, \dots$  the gamma function can be defined by

$$\Gamma(z) = \sum_{n=0}^{\infty} \frac{(-1)^n}{n!} \frac{1}{z+n} + \int_1^{\infty} t^{z-1} e^{-t} dt.$$

Then  $\Gamma(z)$  is a meromorphic function having simple poles at  $0, -1, -2, \dots$ , and

$$\Gamma(z) = \int_0^{\infty} t^{z-1} e^{-t} dt$$

for  $\text{Re}(z) > 0$ . Also  $\Gamma(z) \neq 0$  for all  $z$ . The subroutines CGAMMA and DCGAMA are available for computing  $\Gamma(z)$  and  $\ln \Gamma(z)$  when  $z$  is complex, and the functions GAMMA, GAMMLN, DGAMMA, and DGAMMLN are available for computing  $\Gamma(z)$  and  $\ln \Gamma(z)$  when  $z$  is real. DCGAMA, DGAMMA, and DGAMMLN are double precision procedures.

### CALL CGAMMA(MO, z, w)

MO is an integer,  $z$  a complex number satisfying  $z \neq 0, -1, -2, \dots$ , and  $w$  a complex variable. When CGAMMA is called,  $w$  is assigned the value  $\Gamma(z)$  if  $\text{MO} = 0$  and the value  $\ln \Gamma(z)$  if  $\text{MO} \neq 0$ .

**Error return.** If  $z = 0, -1, -2, \dots$ , or if  $\text{Re}(z) < 0$  and  $\text{Re}(z)$  is too large for  $\ln \Gamma(z)$  to be computed, then  $w$  is assigned the value 0.

**Programming.** CGAMMA calls the functions REXP, SPMPAR, and IPMPAR. CGAMMA was written by A. H. Morris.

### References.

- (1) Kuki, Hirono, "Complex Gamma Function with Error Control," *Comm. ACM* 15 (1972), pp. 262-267.
- (2) Spira, Robert, "Calculation of the Gamma Function by Stirling's Formula," *Math Comp.* 25 (1971), pp. 317-322.

### GAMMA(x)

The argument  $x$  is a real number. If  $\Gamma(x)$  can be computed then GAMMA( $x$ ) is assigned the value  $\Gamma(x)$ . Otherwise, if  $\Gamma(x)$  cannot be computed, then GAMMA( $x$ ) is set to 0.

**Algorithm.** If  $|x| < 15$  then  $x$  is reduced to the interval  $[1, 2]$  by  $\Gamma(a+1) = a\Gamma(a)$ , and a rational minimax approximation is employed. If  $x < -15$  then

$$(1) \quad \Gamma(x) = \frac{(-1)^{n+1} \pi}{\sin(\pi \lambda) |x| \Gamma(|x|)}$$



is applied. Here  $|x| = n + \lambda$  where  $n$  is the largest integer less than  $|x|$ . For  $x \geq 15$

$$(2) \quad \ln \Gamma(x) = (x - 1/2) \ln x - x + \frac{1}{2} \ln(2\pi) + \Delta(x)$$

is computed where  $\Delta(x)$  is a minimax approximation. The function  $\Delta(x)$  is evaluated in single precision, and a double precision value is obtained for  $\ln x$ . This yields a double precision value for  $\ln \Gamma(x)$ . If  $\ln \Gamma(x) = \alpha + \delta$  where  $\alpha$  is the leading portion of  $\ln \Gamma(x)$ , then  $\Gamma(x)$  is set to  $e^\alpha(1 + \delta)$ . This is permissible since  $1 + \delta$  is the portion of the Taylor series expansion for  $e^\delta$  that is significant.

The logarithm  $\ln x$  is evaluated as follows: Let  $n$  be the largest integer less than or equal to  $x$ , and let  $t = (x - n)/(x + n)$ . Then  $x = n(1 + t)/(1 - t)$  so that  $\ln x = \ln n + \ln[(1 + t)/(1 - t)]$ . Also  $0 \leq t < 1/(2n)$ . The function  $\ln [(1 + t)/(1 - t)]$  is computed by a polynomial minimax approximation in single precision, and the value  $\ln n$  is stored in double precision.

**Precision.** If  $0 < x \leq 2$  then GAMMA( $x$ ) is accurate to within 2 units of the 14<sup>th</sup> significant digit. If  $x > 2$  then GAMMA( $x$ ) is accurate to within 3 units of the 14<sup>th</sup> significant digit. Otherwise, GAMMA( $x$ ) is accurate to within 5 units of the 14<sup>th</sup> significant digit.

**Programming.** GAMMA calls the functions GLOG and EXPARG. These functions were written by A. H. Morris. The functions SPMPAR and IPMPAR are also used.

#### GAMLN( $x$ )

GAMLN( $x$ ) =  $\ln \Gamma(x)$  for all positive real  $x$ .

**Algorithm.** See p. 379 and appendix D of the reference.

**Precision.** GAMLN( $x$ ) is accurate to within 2 units of the 14<sup>th</sup> significant digit when GAMLN( $x$ )  $\neq 0$ .

**Reference.** DiDonato, A. R. and Morris, A. H., "Computation of the Incomplete Gamma Function Ratios and Their Inverse," *ACM Trans. Math Software* 12 (1986), pp. 377-393.

**Programming.** GAMLN calls the function GAMLN1. These functions were written by A. H. Morris.

#### CALL DCGAMA(MO, Z, W)

MO is an integer, and  $Z$  and  $W$  are double precision arrays of dimension 2. It is assumed that  $Z(1)$  and  $Z(2)$  are the real and imaginary parts of a complex number  $z$ . If  $MO = 0$  then the double precision value for  $w = \Gamma(z)$  is computed. Otherwise, if  $MO \neq 0$  then the double precision value for  $w = \ln \Gamma(z)$  is computed.  $W(1)$  and  $W(2)$  contain the real and imaginary parts of  $w$ , respectively.

**Error Return.** If  $z = 0, -1, -2, \dots$ , or if  $\text{Re}(z) < 0$  and  $\text{Re}(z)$  is too large for  $\ln \Gamma(z)$  to be computed, then  $W(1)$  and  $W(2)$  are assigned the value 0.

**Programming.** DCGAMA calls the functions DREXP, DPMPAR, and IPMPAR. DCGAMA was written by A. H. Morris.

**References.**

- (1) Kuki, Hirono, "Complex Gamma Function with Error Control," *Comm. ACM* 15 (1972), pp. 262-267.
- (2) Spira, Robert, "Calculation of the Gamma Function by Stirling's Formula," *Math Comp.* 25 (1971), pp. 317-322.

**DGAMMA( $x$ )**

The argument  $x$  is a double precision real number. If  $\Gamma(x)$  can be computed then DGAMMA( $x$ ) is the double precision value for  $\Gamma(x)$ . Otherwise, if  $\Gamma(x)$  cannot be computed, then DGAMMA( $x$ ) is set to 0.

**Remark.** DGAMMA must be declared in calling program to be of type DOUBLE PRECISION.

**Algorithm.** If  $|x| \leq 20$  then  $x$  is reduced to the interval  $[1, 2)$  by  $\Gamma(a+1) = a\Gamma(a)$ , and the Chebyshev expansion in the SLATEC library given by Wayne Fullerton (Los Alamos) is used. If  $x < -20$  then (1) is applied, and if  $x > 20$  then (2) is used.  $\Delta(x)$  is computed by the power series corresponding to the Chebyshev expansion in the SLATEC library given by Wayne Fullerton. The power series was obtained by A. H. Morris.

**Precision.** If  $0 < x \leq 2$  then DGAMMA( $x$ ) is accurate to within 1 unit of the 28<sup>th</sup> significant digit. Otherwise, if  $x > 2$  then DGAMMA( $x$ ) is accurate to within 1 unit of the 25<sup>th</sup> significant digit.

**Programming.** DGAMMA calls the functions DCSEVL, DPDEL, and DXPAR. These functions were written by A. H. Morris. The functions DPMPAR and IPMPAR are also used.

**DGAMLN( $x$ )**

The argument  $x$  is a double precision positive real number. DGAMLN( $x$ ) is the double precision value for  $\ln \Gamma(x)$ .

**Remark.** DGAMLN must be declared in the calling program to be of type DOUBLE PRECISION.

**Algorithm.** If  $.5 \leq x \leq 2.5$  then the Taylor series for  $1/\Gamma(1+a)$  is used, and if  $x \geq 10$  then (2) is applied.  $\Delta(x)$  is computed by the power series corresponding to the Chebyshev expansion in the SLATEC library given by Wayne Fullerton (Los Alamos). The power series was obtained by A. H. Morris.

**Precision.** DGAMLN( $x$ ) is accurate to within 2 units of the 28<sup>th</sup> significant digit when DGAMLN( $x$ )  $\neq$  0.

**Programming.** DGAMLN calls the functions DPDEL, DGMLN1, and DLNREL. These functions were written by A. H. Morris. The function DPMPAR is also used.

## DIGAMMA FUNCTION

For any complex  $z \neq 0, -1, -2, \dots$  the digamma (or psi) function  $\psi(z)$  is defined by

$$\psi(z) = \Gamma'(z)/\Gamma(z)$$

where  $\Gamma(z)$  is the gamma function. For real  $x > 0$ ,  $\psi(x)$  is an increasing function having a zero at the point  $x_0 = 1.4616\ 32144\ 96836$ . The subroutines CPSI and DCPSI are available for computing  $\psi(z)$  when  $z$  is complex, and the functions PSI and DPSI are available for computing  $\psi(z)$  when  $z$  is real. DCPSI and DPSI are double precision procedures.

**CALL CPSI( $z, w$ )**

The argument  $z$  is a complex number satisfying  $z \neq 0, -1, -2, \dots$ , and  $w$  is a complex variable. When CPSI is called,  $w$  is assigned the value  $\psi(z)$ .

**Error Return.** If  $z = 0, -1, -2, \dots$ , or if  $\text{Re}(z) < 0$  and  $\text{Re}(z)$  is too large for  $\psi(z)$  to be computed, then  $w$  is assigned the value 0.

**Algorithm.** If  $z = x + iy$  satisfies  $x \geq 0$  and  $|z| \geq 6$ , then the asymptotic expansion

$$(1) \quad \psi(z) = \ln z - \frac{1}{2z} - \sum_{m=1}^{\infty} \frac{B_{2m}}{2mz^{2m}}$$

is employed. Otherwise, if  $x \geq 0$  then the smallest nonnegative integer  $n$  is found for which  $|z + n| \geq 6$ , and the relation

$$\psi(z) = - \sum_{j=0}^{n-1} \frac{1}{z+j} + \psi(z+n)$$

is applied. When  $x < 0$  then

$$(2) \quad \psi(z) = \psi(1-z) - \pi \cot(\pi z)$$

is also used.

**Programming.** CPSI calls the functions REXP, SPMPAR, and IPMPAR. CPSI was written by A. H. Morris.

**PSI( $x$ )**

The argument  $x$  is a real number. If  $\psi(x)$  can be computed then PSI( $x$ ) is assigned the value  $\psi(x)$ . Otherwise, if  $\psi(x)$  cannot be computed, then PSI( $x$ ) is set to 0.

**Precision.** If  $x > 0$  then PSI( $x$ ) is accurate to within 2 units of the 14<sup>th</sup> significant digit when PSI( $x$ )  $\neq 0$ .

**Programming.** PSI calls the functions SPMPAR and IPMPAR. PSI was written at Argonne National Laboratory for the FUNPACK package of special function subroutines. PSI was modified by A. H. Morris.

**Reference.** Cody, W. J., Strecok, A. J., and Thacher, H. C., "Chebyshev Approximations for the Psi Function," *Math Comp.* 27 (1973), pp. 123-127.

### CALL DCPSI(Z,W)

$Z$  and  $W$  are double precision arrays of dimension 2. It is assumed that  $Z(1)$  and  $Z(2)$  are the real and imaginary parts of a complex number  $z$ . When DCPSI is called the double precision value for  $w = \psi(z)$  is computed.  $W(1)$  and  $W(2)$  contain the real and imaginary parts of  $w$ , respectively.

**Error Return.** If  $z = 0, -1, -2, \dots$ , or if  $\text{Re}(z) < 0$  and  $\text{Re}(z)$  is too large for  $\psi(z)$  to be computed, then  $W(1)$  and  $W(2)$  are assigned the value 0.

**Programming.** DCPSI calls the functions DREXP, DPMPAR, and IPMPAR. DCPSI was written by A. H. Morris.

### DPSI(x)

The argument  $x$  is a double precision real number. If  $\psi(x)$  can be computed then  $\text{DPSI}(x)$  is the double precision value for  $\psi(x)$ . Otherwise, if  $\psi(x)$  cannot be computed, then  $\text{DPSI}(x)$  is set to 0.

**Remark.** DPSI must be declared in the calling program to be of type DOUBLE PRECISION.

**Algorithm.** If  $|x| < 10$  then  $x$  is reduced to the interval  $[1, 2)$  by  $\psi(a+1) = \frac{1}{a} + \psi(a)$ , and the Chebyshev expansion in the SLATEC library given by Wayne Fullerton (Los Alamos) is used when  $|x - x_0| > 2 \cdot 10^{-2}$ . Otherwise, if  $|x - x_0| < 2 \cdot 10^{-2}$  then the Taylor series around the zero  $x_0$  is used. The coefficients for the Taylor series were obtained by A. H. Morris. If  $x < -10$  then (2) is applied, and if  $x \geq 10$  then  $\psi(x) - \ln x$  is computed by the power series corresponding to the Chebyshev expansion in the SLATEC library given by Wayne Fullerton. The power series was obtained by A. H. Morris.

**Precision.** If  $0 < x \leq 1$  or  $x \geq 2$  then  $\text{DPSI}(x)$  is accurate to within 2 units of the 28<sup>th</sup> significant digit. Otherwise, if  $1 < x < 2$  then  $\text{DPSI}(x)$  is accurate to within 5 units of the 28<sup>th</sup> significant digit when  $\text{DPSI}(x) \neq 0$ .

**Programming.** DPSI calls the functions DCSEVL, DPSI0, DPMPAR, and IPMPAR. DPSI was written by A. H. Morris.

## LOGARITHM OF THE BETA FUNCTION

For  $a, b > 0$  the beta function  $B(a, b)$  can be defined by

$$B(a, b) = \int_0^1 t^{a-1} (1-t)^{b-1} dt.$$

From this it follows that  $B(a, b) = \Gamma(a)\Gamma(b)/\Gamma(a+b)$  where  $\Gamma(a)$  is the gamma function. The functions BETALN and DBETLN are available for computing  $\ln B(a, b)$ . DBETLN is a double precision function.

### BETALN( $a, b$ )

$\text{BETALN}(a, b) = \ln B(a, b)$  for  $a, b > 0$ .

**Algorithm.** See pages 19–21 of the reference.

**Precision.**  $\text{BETALN}(a, b)$  is accurate to within 4 units of the 14<sup>th</sup> significant digit when  $a, b \geq 1$  and  $\text{BETALN}(a, b) \neq 0$ . In particular, when  $a, b \geq 15$ ,  $\text{BETALN}(a, b)$  is accurate to within 2 units of the 14<sup>th</sup> significant digit.

**Programming.** BETALN employs the functions ALNREL, ALGDIV, BCORR, GAMLN, GAMLN1, and GSUMLN. These functions were written by A. H. Morris.

**Reference.** DiDonato, A. R. and Morris, A. H., *Significant Digit Computation of the Incomplete Beta Function Ratios*, Report TR 88-365, Naval Surface Warfare Center, Dahlgren, Virginia, 1988.

### DBETLN( $a, b$ )

The arguments  $a$  and  $b$  are positive double precision numbers.  $\text{DBETLN}(a, b)$  is the double precision value for  $\ln B(a, b)$ .

**Remark.** DBETLN must be declared in the calling program to be of type DOUBLE PRECISION.

**Algorithm.** The algorithm for  $\ln B(a, b)$  on pages 19–21 of the reference is used.  $\Delta(x)$  is computed by the power series corresponding to the Chebyshev expansion in the SLATEC library given by Wayne Fullerton (Los Alamos). The power series was obtained by A. H. Morris.

**Programming.** DBETLN employs the functions DLNREL, DLGDIV, DBCORR, DPDEL, DGAMLN, DGMLN1, and DGSMLN. These functions were written by A. H. Morris. The function DPMPAR is also used.

**Reference.** DiDonato, A. R. and Morris, A. H., *Significant Digit Computation of the Incomplete Beta Function Ratios*, Report TR 88-365, Naval Surface Warfare Center, Dahlgren, Virginia, 1988.

## INCOMPLETE GAMMA RATIO FUNCTIONS

For  $a > 0$  and  $x \geq 0$  let  $P(a, x)$  and  $Q(a, x)$  denote the functions defined by

$$P(a, x) = \frac{1}{\Gamma(a)} \int_0^x e^{-t} t^{a-1} dt$$

$$Q(a, x) = \frac{1}{\Gamma(a)} \int_x^\infty e^{-t} t^{a-1} dt.$$

Then  $0 \leq P(a, x) \leq 1$  and  $P(a, x) + Q(a, x) = 1$ . Also,  $P(a, x) \rightarrow 1$  and  $Q(a, x) \rightarrow 0$  for  $x > 0$  when  $a \rightarrow 0$ . Hence, we may define  $P(0, x) = 1$  and  $Q(0, x) = 0$  for  $x > 0$ . The subroutine GRATIO is available for computing  $P(a, x)$  and  $Q(a, x)$ , and the auxiliary function RCOMP is provided for computing  $e^{-x} x^a / \Gamma(a)$ .

### CALL GRATIO( $a, x, P, Q, i$ )

It is assumed that  $a \geq 0$  and  $x \geq 0$ , where  $a$  and  $x$  are not both 0.  $P$  and  $Q$  are variables. GRATIO assigns  $P$  the value  $P(a, x)$  and  $Q$  the value  $Q(a, x)$ . The argument  $i$  may be any integer. This argument specifies the desired accuracy of the results. If  $i = 0$  then the user is requesting as much accuracy as possible (up to 14 significant digits). Otherwise, if  $i = 1$  then accuracy is requested to within 1 unit of the 6<sup>th</sup> significant digit, and if  $i \neq 0, 1$  then the accuracy is requested to within 1 unit of the 3<sup>rd</sup> significant digit.

**Error Return.**  $P$  is assigned the value 2 when  $a$  or  $x$  is negative, when  $a = x = 0$ , or when  $P(a, x)$  and  $Q(a, x)$  are indeterminate.  $P(a, x)$  and  $Q(a, x)$  are indeterminate when  $x \approx a$  and  $a$  is exceedingly large. On the CDC 6000-7000 series computers this occurs when  $|x/a - 1| \leq 10^{-14}$  and  $a \geq 6.6E25$ .

**Programming.** GRATIO calls the functions ERF, ERFC1, REXP, RLOG, GAMMA, GAM1, and SPMPAR. GAMMA employs the functions GLOG, EXPARG, and IPMPAR. GRATIO was written by A. H. Morris.

**Reference.** DiDonato, A. R. and Morris, A. H., "Computation of the Incomplete Gamma Function Ratios and Their Inverse," *ACM Trans. Math Software* 12 (1986), pp. 377-393.

### RCOMP( $a, x$ )

$RCOMP(a, x) = e^{-x} x^a / \Gamma(a)$  for  $a > 0$  and  $x > 0$ .

**Algorithm.** See page 378 of the reference.

**Programming.** RCOMP employs the functions EXPARG, GAMMA, GAM1, GLOG, and RLOG. These functions were written by A. H. Morris. The functions SPMPAR and IPMPAR are also used.

**Reference.** DiDonato, A. R. and Morris, A. H., "Computation of the Incomplete Gamma Function Ratios and their Inverse," *ACM Trans. Math Software* 12 (1986), pp. 337-393.

## INVERSE INCOMPLETE GAMMA RATIO FUNCTION

For  $a > 0$  and  $x \geq 0$  let  $P(a, x)$  and  $Q(a, x)$  denote the incomplete gamma ratio functions defined by

$$P(a, x) = \frac{1}{\Gamma(a)} \int_0^x e^{-t} t^{a-1} dt$$

$$Q(a, x) = \frac{1}{\Gamma(a)} \int_x^\infty e^{-t} t^{a-1} dt.$$

Then  $0 \leq P(a, x) \leq 1$  and  $P(a, x) + Q(a, x) = 1$ . If we are given  $a, p$ , and  $q$  where  $a > 0, 0 \leq p \leq 1$ , and  $p + q = 1$ , then the subroutine GAMINV is available for obtaining the value  $x \geq 0$  for which  $P(a, x) = p$  and  $Q(a, x) = q$ .

**CALL GAMINV( $a, X, x_0, p, q, \text{IND}$ )**

$X$  is a variable. If  $p = 0$  then  $X$  is assigned the value 0, and if  $q = 0$  then  $X$  is set to the largest floating point number available. Otherwise, GAMINV attempts to obtain a solution  $x$  for  $P(a, x) = p$  and  $Q(a, x) = q$  that is correct to at least 10 significant digits.<sup>1</sup> If the routine is successful then the solution is stored in  $X$ . The solution is normally obtained by Schroder iteration. The argument  $x_0$  is an optional initial approximation for  $x$ . If the user does not wish to supply an initial approximation then set  $x_0 \leq 0$ .

IND is a variable that reports the status of the results. When GAMINV terminates, IND has one of the following values:

- IND = 0 The solution was obtained. Iteration was not used.
- IND  $\geq$  1 The solution was obtained. IND iterations were performed.
- IND = -2 (Input error)  $a \leq 0$ .
- IND = -3 No solution was obtained. The ratio  $Q/a$  is too large.
- IND = -4 (Input error)  $p + q \neq 1$ .
- IND = -6 20 iterations were performed. The most recent value obtained for  $x$  is stored in  $X$ . This cannot occur if  $x_0 \leq 0$ .
- IND = -7 Iteration failed. No value is given for  $x$ . This may occur when  $x \approx 0$ .
- IND = -8 A value for  $x$  is stored in  $X$ , but the routine is not certain of its accuracy. Iteration cannot be performed in this case. If  $x_0 \leq 0$  then this can occur only when  $p \approx 0$  or  $q \approx 0$ . If  $x_0 > 0$  then this can occur when  $a \approx x$  and  $a$  is exceedingly large (say  $a \geq 10^{20}$ ).

**Remark.** If  $x_0 \leq 0$  then 3 or fewer iterations are required.

**Programming.** GAMINV employs the routine GRATIO and functions ERF, ERFC1, REXP, RLOG, ALNREL, GAMMA, GAM1, GAMLN, GAMLN1, RCOMP, and SPMPAR.

<sup>1</sup>If a  $k$  digit floating point arithmetic is being used where  $k < 10$ , then the routine attempts to obtain a solution that is correct to machine accuracy.



GAMMA uses the functions GLOG, EXPARG, and IPMPAR. GAMINV was written by A. H. Morris

**Reference.** DiDonato, A. R. and Morris, A. H., "Computation of the Incomplete Gamma Function Ratios and Their Inverse," *ACM Trans. Math Software* 12 (1986), pp. 377-393.

## INCOMPLETE BETA FUNCTION

For  $a, b > 0$  and  $0 \leq x \leq 1$  the incomplete beta function is defined by

$$I_x(a, b) = \frac{1}{B(a, b)} \int_0^x t^{a-1} (1-t)^{b-1} dt$$

where  $B(a, b)$  is the beta function. Then we note that  $0 \leq I_x(a, b) \leq 1$  and

$$\begin{aligned} \lim_{a \rightarrow 0} I_x(a, b) &= 1 \text{ for } x \neq 0 \\ \lim_{b \rightarrow 0} I_x(a, b) &= 0 \text{ for } x \neq 1. \end{aligned}$$

These limits permit  $I_x(a, b)$  to be defined to be 1 when  $a = 0$  and  $b \neq 0, x \neq 0$ , and for  $I_x(a, b)$  to be defined to be 0 when  $b = 0$  and  $a \neq 0, x \neq 1$ . The subroutine BRATIO is available for computing  $I_x(a, b)$  for arbitrary  $a, b \geq 0$ , and the subroutine ISUBX is available for computing  $I_x(a, b)$  for the highly specialized case when  $a$  and  $b$  are integers or half-integers. Also, the auxiliary function BRCOMP is provided for computing  $x^a y^b / B(a, b)$  when  $0 < x < 1$  and  $y = 1 - x$ .

**CALL BRATIO( $a, b, x, y, W, W1, IERR$ )**

It is assumed that  $a \geq 0, b \geq 0, 0 \leq x \leq 1$ , and  $y = 1 - x$ .  $W, W1$ , and  $IERR$  are variables. If no input errors are detected then  $IERR$  is set to 0,  $W$  is assigned the value  $I_x(a, b)$ , and  $W1$  is assigned the value  $1 - I_x(a, b)$ .

**Error Return.** When an input error is detected, then  $W$  and  $W1$  are assigned the value 0 and  $IERR$  is set to one of the following values:

- $IERR = 1$  if  $a < 0$  or  $b < 0$
- $IERR = 2$  if  $a = b = 0$
- $IERR = 3$  if  $x < 0$  or  $x > 1$
- $IERR = 4$  if  $y < 0$  or  $y > 1$
- $IERR = 5$  if  $x + y \neq 1$
- $IERR = 6$  if  $x = a = 0$
- $IERR = 7$  if  $y = b = 0$

**Programming.** BRATIO employs the subroutines BGRAT and GRAT1, and the functions ALGDIV, ALNREL, BASYM, BCORR, BETALN, BFRAC, BPSER, BRCOMP, BRCMP1, BUP, ERF, ERFC1, GAMLN, GAMLN1, GAM1, GSUMLN, ESUM, EXPARG, REXP, and RLOG1. These subroutines and functions were written by A. H. Morris. The functions SPMPAR and IPMPAR are also used.

**Reference.** DiDonato, A. R. and Morris, A. H., *Significant Digit Computation of the Incomplete Beta Function Ratios*, Report TR 88-365, Naval Surface Warfare Center, Dahlgren, Virginia, 1988.

### CALL ISUBX( $a, b, x, W, IERR, EPS$ )

It is assumed that  $a, b$ , and  $x$  satisfy the following restrictions:

- (1)  $a > 0, b > 0$ , and  $x \geq 0$
- (2)  $a \geq 1/2, 1/2 \leq b \leq 70$ , and  $x \leq 1$
- (3)  $a$  and  $b$  are integers or half-integers

EPS specifies the (absolute) accuracy that is desired.  $W$  is a real variable and IERR an integer variable. When ISUBX is called, if there are no input errors then  $W$  is assigned the value  $I_x(a, b)$  and IERR is assigned the value 1.

**Error Return.** If an error is detected then IERR is assigned one of the following values:

- IERR = 2 if restrictions (1) are violated.
- IERR = 3 if restrictions (2) are violated or  $a$  is too large.
- IERR = 4 if restrictions (3) are violated.

Also  $W$  is assigned the value 0.

**Remarks.** ISUBX was designed for a maximum precision  $EPS = 10^{-10}$

**Programming.** ISUBX employs the functions ALGDIV, ALNREL, BLND, IPMPAR, and LOGAM. ISUBX was written by A. H. Morris.

**Reference.** DiDonato, A. R. and Jarnagin, M. P., "The Efficient Calculation of the Incomplete Beta-function Ratio for Half-Integer Values of the Parameter  $a, b$ ," *Math Comp.* 21 (1967), pp. 652-662.

### BRCOMP( $a, b, x, y$ )

$BRCOMP(a, b, x, y) = x^a y^b / B(a, b)$  for  $a, b > 0$  and  $x, y > 0$  where  $x + y = 1$ .

**Algorithm.** See pages 19-21 of the reference.

**Programming.** BRCOMP employs the functions ALGDIV, ALNREL, BCORR, BETALN, GAM1, GAMLN, GAMLN1, GSUMLN, and RLOG1. These functions were written by A. H. Morris.

**Reference.** DiDonato, A. R. and Morris, A. H., *Significant Digit Computation of the Incomplete Beta Function Ratios*, Report TR 88-365, Naval Surface Warfare Center, Dahlgren, Virginia, 1988.

## BESSEL FUNCTION $J_\nu(z)$

If  $\nu$  is complex then  $J_\nu(z)$  is defined by

$$J_\nu(z) = \sum_{k=0}^{\infty} \frac{(-1)^k (z/2)^{\nu+2k}}{k! \Gamma(\nu+k+1)}$$

for any  $z \neq 0$  in the complex plane cut along the negative real axis.  $J_\nu(z)$  is analytic in the region  $|\arg(z)| < \pi$ , and  $J_\nu(z)$  is an entire function of  $\nu$  for any fixed  $z$ . If  $\nu$  is an integer then  $J_\nu(z)$  is also defined at 0 and is an entire function of  $z$ . The following subroutines are available for computing  $J_\nu(z)$ .

### CALL CBSSLJ( $z, \nu, w$ )

The arguments  $z$  and  $\nu$  are complex numbers and  $w$  is a complex variable. It is assumed that  $|\arg(z)| < \pi$ . When CBSSLJ is called,  $w$  is assigned the value  $J_\nu(z)$ .

**Precision.** CBSSLJ is accurate to within  $4 \cdot 10^{-13}$  for real  $0 < z \leq 35$  and  $0 \leq \nu \leq 1$ .

**Note.** CBSSLJ employs the subroutine CGAMMA.

**Programmer.** A. V. Hershey

**Reference.** Hershey, A. V., *Computation of Special Functions*, Report TR-3788, Naval Surface Weapons Center, Dahlgren, Virginia, 1978.

### CALL BSSLJ( $z, n, w$ )

The argument  $z$  is a complex number,  $n$  is an integer, and  $w$  is a complex variable. When BSSLJ is called,  $w$  is assigned the value  $J_n(z)$ .

**Precision.** BSSLJ is accurate to within  $5 \cdot 10^{-14}$  for real  $0 < z \leq 35$  and  $n = 0, 1$ .

**Programmer.** A. V. Hershey

**Reference.** Hershey, A. V., *Computation of Special Functions*, Report TR-3788, Naval Surface Weapons Center, Dahlgren, Virginia, 1978.

### CALL BESJ( $x, \alpha, n, W, k$ )

The arguments  $x$  and  $\alpha$  are nonnegative real numbers,  $n$  is a positive integer, and  $W$  is an array of dimension  $n$  or larger. When BESJ is called  $J_{\alpha+i-1}(x)$  is computed and stored in  $W(i)$  for  $i = 1, \dots, n$ .

The argument  $k$  is an integer variable that is set by the routine. If all  $J_{\alpha+i-1}(x)$  are computed then  $k$  is set to 0. Otherwise,  $k$  is assigned one of the following values:

$k = -1$  The argument  $x$  is negative.

$k = -2$  The argument  $\alpha$  is negative.

$k = -3$  The requirement  $n \geq 1$  is violated.

$k > 0$  The last  $k$  components of  $W$  have been set to 0 because of underflow.

**Precision.** For  $0 < x \leq 35$  and  $0 \leq \alpha \leq 1$ , BESJ is accurate to within  $8 \cdot 10^{-13}$ .

**Programming.** BESJ calls the subroutines ASJY and JAIRY, and the functions GAMLN, SPMPAR, and IPMPAR. The subroutines were written by Donald E. Amos, Sharon L. Daniel, and M. Katherine Weston (Sandia Laboratories).

#### References.

- (1) Amos, D.E., Daniel, S. L., and Weston, M. K., *CDC 6600 Subroutines for Bessel Functions  $J_\nu(x)$ ,  $x \geq 0$ ,  $\nu \geq 0$  and Airy Functions  $A_i(x)$ ,  $A_i'(x)$ ,  $-\infty < x < \infty$* . Report SAND 75-0147, Sandia Laboratories, Albuquerque, New Mexico, 1975.
- (2) \_\_\_\_\_, "CDC 6600 Subroutines IBESS and JBESS for Bessel Functions  $I_\nu(x)$  and  $J_\nu(x)$ ,  $x \geq 0$ ,  $\nu \geq 0$ ," *ACM Trans. Math Software* **3** (1977), pp. 76-92.

## BESSEL FUNCTION $Y_\nu(z)$

If  $\nu$  is any complex number not an integer, then  $Y_\nu(z)$  can be defined by

$$Y_\nu(z) = \frac{J_\nu(z) \cos \nu\pi - J_{-\nu}(z)}{\sin \nu\pi}$$

for any  $z \neq 0$  in the complex plane cut along the negative real axis. For any integer  $n$  we can also define  $Y_n(z) = \lim_{\nu \rightarrow n} Y_\nu(z)$ . Then for any complex  $\nu$ ,  $Y_\nu(z)$  is analytic in the region  $|\arg(z)| < \pi$ . Also,  $Y_\nu(z)$  is an entire function of  $\nu$  for any fixed  $z$ . The following subroutine is available for computing  $Y_\nu(z)$  when  $\nu$  is an integer.

**CALL BSSLY( $z, n, w$ )**

The argument  $z$  is a complex number,  $n$  is an integer, and  $w$  is a complex variable. It is assumed that  $|\arg(z)| < \pi$ . When BSSLY is called,  $w$  is assigned the value  $Y_n(z)$ .

**Precision.** If  $.005 \leq x \leq .785$  then  $Y_0(x)$  and  $Y_1(x)$  are accurate to within 3 units of the 14<sup>th</sup> significant digit. Otherwise, if  $x > .785$  then  $Y_0(x)$  and  $Y_1(x)$  are accurate to within  $4 \cdot 10^{-14}$ .

**Programmer.** A. V. Hershey

**Reference.** Hershey, A. V., *Computation of Special Functions*, Report TR-3788, Naval Surface Weapons Center, Dahlgren, Virginia, 1978.

## MODIFIED BESSEL FUNCTION $I_\nu(z)$

If  $\nu$  is complex then  $I_\nu(z)$  is defined by

$$I_\nu(z) = \sum_{k=0}^{\infty} \frac{(z/2)^{\nu+2k}}{k! \Gamma(\nu+k+1)}$$

for any  $z \neq 0$  in the complex plane cut along the negative real axis.  $I_\nu(z)$  is analytic in the region  $|\arg(z)| < \pi$ , and  $I_\nu(z)$  is an entire function of  $\nu$  for any fixed  $z$ . If  $\nu$  is an integer then  $I_\nu(z)$  is also defined at 0 and is an entire function of  $z$ . The following subroutines are available for computing  $I_\nu(z)$ .

### CALL BSSLI(MO, $z$ , $n$ , $w$ )

MO is an integer,  $z$  a complex number,  $n$  an integer, and  $w$  a complex variable. If MO  $\neq 0$  then it is assumed that  $|\arg(z)| < \pi$ . When BSSLI is called,  $w$  is assigned the value  $I_n(z)$  if MO = 0 and the value  $e^{-z} I_n(z)$  if MO  $\neq 0$ .

**Precision.** BSSLI is accurate to within 5 units of the 13<sup>th</sup> significant digit for real  $0 < z \leq 35$  and  $n = 0, 1, \dots, 40$ .

**Programmer.** Allen V. Hershey

**Reference.** Hershey, A. V., *Computation of Special Functions*, Report TR-3788, Naval Surface Weapons Center, Dahlgren, Virginia, 1978.

### CALL BESI( $x$ , $\alpha$ , MO, $n$ , $W$ , $k$ )

MO may be 1 or 2. The arguments  $x$  and  $\alpha$  are nonnegative real numbers,  $n$  is a positive integer, and  $W$  is an array of dimension  $n$  or larger. When BESI is called, if MO = 1 then  $I_{\alpha+i-1}(x)$  is computed and stored in  $W(i)$  for  $i = 1, \dots, n$ . Otherwise, if MO = 2 then  $e^{-x} I_{\alpha+i-1}(x)$  is computed and stored in  $W(i)$ .

The argument  $k$  is an integer variable that is set by the routine. If all  $I_{\alpha+i-1}(x)$  or  $e^{-x} I_{\alpha+i-1}(x)$  are computed then  $k$  is set to 0. Otherwise,  $k$  is assigned one of the following values:

- $k = -1$  The argument  $x$  is negative.
- $k = -2$  The argument  $\alpha$  is negative.
- $k = -3$  The requirement  $n \geq 1$  is violated.
- $k = -4$  MO is not 1 or 2.
- $k = -5$  The argument  $x$  is too large for MO = 1.
- $k > 0$  The last  $k$  components of  $W$  have been set to 0 because of underflow.

**Precision.** For  $0 < x \leq 35$  and  $0 \leq \alpha \leq 1$ , or  $0 < x \leq 35$  and  $\alpha = 1, 2, \dots, 40$ ,  $I_\alpha(x)$  is accurate to within 2 units of the 12<sup>th</sup> significant digit.

**Programming.** BESI calls the subroutine ASIK and the functions GAMLN, SPMPAR, and IPMPAR. BESI and ASIK were written by D. E. Amos and S. L. Daniel (Sandia Laboratories).

**References.**

- (1) Amos, D. E. and Daniel, S. L., *A CDC 6600 Subroutine for Bessel Functions  $I_\nu(x)$ ,  $\nu \geq 0, x \geq 0$* . Report SAND 75-0152, Sandia Laboratories, Albuquerque, New Mexico, 1975.
- (2) Amos, D. E., Daniel, S. L., and Weston, M. K., "CDC 6600 Subroutines IBESS and JBESS for Bessel Functions  $I_\nu(x)$  and  $J_\nu(x)$ ,  $x \geq 0, \nu \geq 0$ ," *ACM Trans. Math Software* **3** (1977), pp. 76-92.



## MODIFIED BESSEL FUNCTION $K_\nu(z)$

If  $\nu$  is any complex number not an integer, then  $K_\nu(z)$  is defined by

$$K_\nu(z) = \frac{\pi}{2} \frac{I_{-\nu}(z) - I_\nu(z)}{\sin \nu\pi}$$

for any  $z \neq 0$  in the complex plane cut along the negative real axis. For any integer  $n$  we can also define  $K_n(z) = \lim_{\nu \rightarrow n} K_\nu(z)$ . Then for any complex  $\nu$ ,  $K_\nu(z)$  is analytic in the region  $|\arg(z)| < \pi$ . Also,  $K_\nu(z)$  is an entire function of  $\nu$  for any fixed  $z$ . The following subroutines are available for computing  $K_\nu(z)$ .

### CALL CBSSLK( $z, r, w$ )

The argument  $z$  is a complex number,  $r$  is a real number, and  $w$  is a complex variable. It is assumed that  $|\arg(z)| < \pi$ . When CBSSLK is called,  $w$  is assigned the value  $K_r(z)$ .

**Programmer.** Allen V. Hershey

**Reference.** Hershey, A. V., *Approximations of Functions by Sets of Poles*, Report TR-2564, Naval Weapons Laboratory, Dahlgren, Virginia, 1971.

### CALL BSSLK( $MO, z, n, w$ )

$MO$  is an integer,  $z$  a complex number,  $n$  an integer, and  $w$  a complex variable. It is assumed that  $|\arg(z)| < \pi$ . When BSSLK is called,  $w$  is assigned the value  $K_n(z)$  if  $MO = 0$  and the value  $e^z K_n(z)$  if  $MO \neq 0$ .

**Precision.** BSSLK is accurate to within 6 units of the 14<sup>th</sup> significant digit for real  $z$  and  $n = 0, 1$ .

**Programmer.** Allen V. Hershey

**Reference.** Hershey, A. V., *Computation of Special Functions*, Report TR-3788, Naval Surface Weapons Center, Dahlgren, Virginia, 1978.

## AIRY FUNCTIONS

For any series  $w = \sum_{n=0}^{\infty} a_n z^n$  satisfying the differential equation  $w'' = zw$ , it follows that  $w = a_0 f(z) + a_1 g(z)$  where

$$f(z) = 1 + \sum_{n \geq 1} \frac{1 \cdot 4 \cdots (3n-2)}{(3n)!} z^{3n}$$

$$g(z) = 1 + \sum_{n \geq 1} \frac{2 \cdot 5 \cdots (3n-1)}{(3n+1)!} z^{3n+1}.$$

In particular, the Airy functions  $Ai(z)$  and  $Bi(z)$  are independent solutions of  $w'' = zw$  where

$$Ai(z) = c_1 f(z) - c_2 g(z)$$

$$Bi(z) = \sqrt{3} [c_1 f(z) + c_2 g(z)]$$

for  $c_1 = 3^{-2/3}/\Gamma(2/3)$  and  $c_2 = 3^{-1/3}/\Gamma(1/3)$ .  $Ai(z)$  and  $Bi(z)$  are entire functions.

The subroutines CAI and CBI are available for computing  $Ai(z)$  and  $Bi(z)$  when  $z$  is complex, and the functions AI, AIE, BI, and BIE are available for computing  $Ai(z)$  and  $Bi(z)$  when  $z$  is real. CAI and CBI also provide the derivatives  $Ai'(z)$  and  $Bi'(z)$  of  $Ai(z)$  and  $Bi(z)$ .

### CALL CAI(IND,z,w,w',IERR)

IND is an integer,  $z$  a complex number, and  $w$  and  $w'$  complex variables. When CAI is called,  $w$  is assigned the value  $Ai(z)$  and  $w'$  the value  $Ai'(z)$  when  $IND = 0$ . Otherwise, if  $IND \neq 0$  then  $w = e^{\zeta} Ai(z)$  and  $w' = e^{\zeta} Ai'(z)$  where  $\zeta = \frac{2}{3}z^{3/2}$ .

IERR is a variable that is set by the routine. When CAI terminates, IERR has one of the following values:

IERR = 0 The desired values were obtained.

IERR = 1  $z$  is too large for the desired values to be computed. In this case  $w$  and  $w'$  are assigned the value 0.

**Precision.** For  $IND \neq 0$  the real and imaginary parts of  $w$  and  $w'$  are accurate to within 2 units of the 12<sup>th</sup> significant digit except near points where they vanish.

**Programming.** CAI employs the subroutines AIRM, AIL, AIA, JA, JMC, BJM, KA, KML, IMC, and BIM. These routines were written by Andrew H. van Tuyl (NSWC) and modified by A. H. Morris. The subroutines CAPO and CREC, and functions CPABS, EXPARG, IPMPAR, and SPMPAR are also used.

### CALL CBI(IND,z,w,w',IERR)

IND is an integer,  $z$  a complex number, and  $w$  and  $w'$  complex variables. When CBI is called,  $w$  is assigned the value  $B_i(z)$  and  $w'$  the value  $B_i'(z)$  when  $IND = 0$ . Otherwise, if  $IND \neq 0$  then

$$w = \begin{cases} e^{-\zeta} B_i(z) & \text{if } |\arg(z)| \leq \pi/3 \\ e^{\zeta} B_i(z) & \text{otherwise} \end{cases}$$

$$w' = \begin{cases} e^{-\zeta} B_i'(z) & \text{if } |\arg(z)| \leq \pi/3 \\ e^{\zeta} B_i'(z) & \text{otherwise} \end{cases}$$

where  $\zeta = \frac{2}{3}z^{3/2}$ .

IERR is a variable that is set by the routine. When CBI terminates, IERR has one of the following values:

IERR = 0 The desired values were obtained.

IERR = 1  $z$  is too large for the desired values to be computed. In this case  $w$  and  $w'$  are assigned the value 0.

**Precision.** For  $IND \neq 0$  the real and imaginary parts of  $w$  and  $w'$  are accurate to within 2 units of the 12<sup>th</sup> significant digit except near points where they vanish.

**Programming.** CBI employs the subroutines AIRM, BII, BIA, IA, IMC, BIM, JA, JMC, and BJM. These routines were written by Andrew H. van Tuyl (NSWC) and modified by A. H. Morris. The subroutine CREC and functions CPABS, EXPARG, IPMPAR, and SPMPAR are also used.

**AI(x)**

$AI(x) = Ai(x)$  for real  $x$ .

**Algorithm.** Rational minimax approximations are used. If  $x < -1$  then  $R$  and  $\theta$  are computed where  $Ai(x) = R \sin(\pi/4 + \theta)$ .

**Precision.** For  $x \geq -1$ ,  $AI(x)$  is accurate to within 2 units of the 12<sup>th</sup> significant digit when  $AI(x) \neq 0$ .

**Programming.** AI calls the subroutine AIMP and function EXPARG. These subprograms were written by A. H. Morris. The function IPMPAR is also used.

**AIE(x)**

If  $x \geq 0$  then  $AIE(x) = e^{\zeta} Ai(x)$  where  $\zeta = \frac{2}{3}x^{3/2}$ . Otherwise, if  $x < 0$  then  $AIE(x) = Ai(x)$ .

**Algorithm.** Rational minimax approximations are used. If  $x < -1$  then  $R$  and  $\theta$  are computed where  $Ai(x) = R \sin(\pi/4 + \theta)$ .

**Precision.** For  $x \geq -1$ ,  $AIE(x)$  is accurate to within 2 units of the 14<sup>th</sup> significant digit.

**Programming.** AIE calls the subroutine AIMP. AIE and AIMP were written by A. H. Morris.

### **BI( $x$ )**

$BI(x) = Bi(x)$  for real  $x$ . If  $x$  is a positive value for which  $Bi(x)$  is too large to be computed, then  $BI(x)$  is assigned the value 0.

**Algorithm.** Rational minimax approximations are used. If  $x < -1$  then  $R$  and  $\theta$  are computed where  $Bi(x) = R \cos(\pi/4 + \theta)$ .

**Precision.** For  $x \geq -1$ ,  $BI(x)$  is accurate to within 2 units of the 12<sup>th</sup> significant digit when  $BI(x) \neq 0$ .

**Programming.** BI calls the subroutine AIMP and function EXPARG. These subprograms were written by A. H. Morris. The function IPMPAR is also used.

### **BIE( $x$ )**

If  $x \geq 0$  then  $BIE(x) = e^{-\zeta} Bi(x)$  where  $\zeta = \frac{2}{3}x^{3/2}$ . Otherwise, if  $x < 0$  then  $BIE(x) = Bi(x)$ .

**Algorithm.** Rational minimax approximations are used. If  $x < -1$  then  $R$  and  $\theta$  are computed where  $Bi(x) = R \cos(\pi/4 + \theta)$ .

**Precision.** For  $x \geq -1$ ,  $BIE(x)$  is accurate to within 2 units of the 14<sup>th</sup> significant digit.

**Programming.** BIE calls the subroutine AIMP. BIE and AIMP were written by A. H. Morris.

## COMPLETE COMPLEX ELLIPTIC INTEGRALS OF THE FIRST AND SECOND KINDS

If  $k$  is complex then the complete elliptic integrals of the first and second kinds can be defined by

$$K(k) = \int_0^{\pi/2} (1 - k^2 \sin^2 t)^{-1/2} dt$$

$$E(k) = \int_0^{\pi/2} (1 - k^2 \sin^2 t)^{1/2} dt$$

for  $|\arg(1 - k^2)| < \pi$ .  $K(k)$  and  $E(k)$  can be extended to  $-\pi \leq \arg(1 - k^2) < \pi$ . For  $|k| < 1$

$$(1) \quad K(k) = \frac{\pi}{2} \sum_{n \geq 0} c_n k^{2n}$$

$$E(k) = -\frac{\pi}{2} \sum_{n \geq 0} c_n \frac{k^{2n}}{2n - 1}$$

where  $c_n = \left[ \frac{(2n)!}{4^n (n!)^2} \right]^2$ . Also, if  $\ell^2 = 1 - k^2$  where  $|\ell| \leq 1$  and  $-\pi \leq \arg(\ell^2) < \pi$ , then

$$(2) \quad K(k) = \frac{1}{\pi} K(\ell) \ln \frac{16}{\ell^2} - \sum_{n \geq 1} c_n \sum_{m=1}^n \frac{\ell^{2n}}{m(2m-1)}$$

$$E(k) = \frac{1}{\pi} [K(\ell) - E(\ell)] \ln \frac{16}{\ell^2} - \sum_{n \geq 1} c_n \frac{2n}{2n-1} \sum_{m=1}^n \frac{\ell^{2n}}{m(2m-1)} + \sum_{n \geq 0} c_n \frac{\ell^{2n}}{(2n-1)^2}.$$

The function CK is available for computing  $K(k)$ , and the subroutine CKE for computing  $K(k)$  and  $E(k)$ .

**CK( $k, \ell$ )**

$CK(k, \ell) = K(k)$  for any complex  $k$  and  $\ell$  where  $k^2 + \ell^2 = 1$  and  $\ell \neq 0$ . CK is a complex valued function which must be declared in the calling program to be of type COMPLEX.

**Error Return.**  $CK(k, \ell) = 0$  if  $\ell = 0$  or  $k^2 + \ell^2 \neq 1$ .

**Remarks.**

- (1)  $CK(k, \ell)$  may underflow, yielding the value 0, when  $|k|$  is sufficiently large.
- (2) CK and the subroutine CKE employ the same algorithm for  $K(k)$ .

**Precision.** If  $k$  is real and  $|k| < 1$  then the relative error of CK is less than  $10^{-13}$ . Also, if  $k$  is purely imaginary then the relative error is less than  $10^{-13}$ .  $K(k)$  is real-valued for only these values of  $k$ . Otherwise, let  $\varepsilon_k = 10^{-12}$  if  $|k| < 0.8$ ,  $\varepsilon_k = 2 \cdot 10^{-13}$  if  $0.8 \leq |k| < 2$ , and  $\varepsilon_k = 10^{-13}$  if  $|k| \geq 2$ . Then the relative errors of the real and imaginary parts of CK are less than  $\varepsilon_k$  except when underflow occurs,  $|k| < 1$  and  $|\arg(\pm k)| < 10^{-287}$ , or  $|k| < 10^{15}$  and  $|\pi/2 - \arg(\pm k)| < 10^{-280}$ . In the latter two cases the relative error of the real part of CK is less than  $\varepsilon_k$ , but all relative accuracy for the imaginary part may be lost.

**Programming.** CK calls the subroutine KL and functions ALNREL, CFLECT, KM, and SPMPAR. CK, KL, and KM were written by Andrew H. van Tuyl (NSWC) and modified by A. H. Morris.

**CALL CKE( $k, \ell, K, E, IERR$ )**

The arguments  $k$  and  $\ell$  are complex numbers where  $k^2 + \ell^2 = 1$  and  $\ell \neq 0$ ,  $K$  and  $E$  are complex variables, and  $IERR$  an integer variable. When CKE is called, if no errors are detected then  $IERR$  is set to 0,  $K$  is assigned the value  $K(k)$ , and  $E$  is assigned the value  $E(k)$ .

**Error Return.**  $IERR = 1$  if  $\ell = 0$  and  $IERR = 2$  if  $k^2 + \ell^2 \neq 1$ . In these cases,  $K$  and  $E$  are not defined.

**Algorithm.** For  $k = 0$  or  $-\pi/2 < \arg(k) \leq \pi/2$ , formulae (2) are used if  $|\ell| \leq .55$ , (1) are used if  $|\ell| > .55$  and  $|k| \leq .55$ , and approximations of the form

$$(3) \quad K(k) = \ell \sum_{n=1}^N \frac{a_n}{\ell^2 + b_n^2 k^2}$$

$$E(k) = \ell \sum_{n=1}^N a_n \left[ 1 + \frac{b_n k}{\ell} \tan^{-1} \frac{b_n k}{\ell} \right]$$

are used if  $|\ell| > .55$  and  $.55 \leq |k| \leq 1$ . (3) are obtained from integral representations for  $K(k)$  and  $E(k)$  by numerical quadrature. If  $|\ell| > .55$ ,  $|k| > 1$ , and  $|k| \leq |\ell|$  then

$$(4) \quad K(k) = \ell_1 K(k_1) \quad k_1 = \pm i k / \ell$$

$$E(k) = E(k_1) / \ell_1 \quad \ell_1 = 1 / \ell$$

are applied where the sign in  $k_1$  is selected so that  $-\pi/2 < \arg(k_1) \leq \pi/2$ . Otherwise, if  $|\ell| > .55$ ,  $|k| > 1$ , and  $|k| > |\ell|$  let  $k_1 = 1/k$  and  $\ell_1 = \pm i \ell / k$  where the sign is selected so that  $-\pi/2 < \arg(\ell_1) \leq \pi/2$ . Then

$$(5) \quad K(k) = k_1 [K(k_1) + i s K(\ell_1)]$$

$$E(k) = \frac{1}{k_1} [E(k_1) - \ell_1^2 K(k_1) - i s (E(\ell_1) - k_1^2 K(\ell_1))]$$

are applied where  $s = 1$  if  $\text{Im}(k) \geq 0$  and  $s = -1$  if  $\text{Im}(k) < 0$ . If  $\arg(k) > \pi/2$  or  $\arg(k) \leq -\pi/2$ , then  $K(k) = K(-k)$  and  $E(k) = E(-k)$  are applied.

**Precision.** If  $k$  is real and  $|k| < 1$ , or  $k$  is purely imaginary, then the relative error of  $E$  is less than  $10^{-13}$ .  $E(k)$  is real-valued for only these values of  $k$ . Otherwise, let  $\varepsilon_k = 10^{-12}$  if  $|k| < 2$  and  $\varepsilon_k = 10^{-13}$  if  $|k| \geq 2$ . Then the relative errors of the real and imaginary parts of  $E$  are less than  $\varepsilon_k$  except when underflow occurs,  $|k| < 1$  and  $|\arg(\pm k)| < 10^{-280}$ , or  $|k| < 10^{15}$  and  $|\pi/2 - \arg(\pm k)| < 10^{-280}$ . In the latter two cases the relative error of the real part of  $E$  is less than  $\varepsilon_k$ , but all relative accuracy for the imaginary part may be lost.

**Programming.** CKE calls the subroutines EKL and EKM, and the functions ALNREL, ATN, CFLECT, and SPMPAR. CKE was written by Andrew H. van Tuyl (NSWC).

## REAL ELLIPTIC INTEGRALS OF THE FIRST AND SECOND KINDS

If  $0 \leq \phi \leq \pi/2$ , then the elliptic integrals of the first and second kinds are defined by

$$F(\phi, k) = \int_0^\phi (1 - k^2 \sin^2 t)^{-1/2} dt$$

$$E(\phi, k) = \int_0^\phi (1 - k^2 \sin^2 t)^{1/2} dt$$

for any real  $k$  where  $k^2 \leq 1$  and  $1 - k^2 \sin^2 \phi \neq 0$ . Alternatively, we may consider

$$R_F(a, b, c) = \frac{1}{2} \int_0^\infty [(t+a)(t+b)(t+c)]^{-1/2} dt$$

where  $a, b, c$  are nonnegative and at most one of them is 0, and

$$R_D(a, b, c) = \frac{3}{2} \int_0^\infty (t+a)^{-1/2} (t+b)^{-1/2} (t+c)^{-3/2} dt$$

where  $a$  and  $b$  are nonnegative such that  $a+b > 0$ , and  $c$  is positive. If  $a \leq b \leq c$  and  $a < c$  then

$$R_F(a, b, c) = \frac{c^{-1/2}}{\sin \phi} F(\phi, k)$$

$$R_D(a, b, c) = \frac{3c^{-3/2}}{k^2 \sin^3 \phi} \left[ F(\phi, k) - E(\phi, k) \right]$$

where  $\cos^2 \phi = a/c$  and  $k^2 = (c-b)/(c-a)$ . If  $\phi = \pi/2$  then the integrals  $F(\phi, k)$  and  $E(\phi, k)$  are said to be *complete*. Otherwise, if  $\phi < \pi/2$  then the integrals are said to be *incomplete*. The subroutines ELLPI, RFVAL, RDVAL, DELLPI, DRFVAL, and DRDVAL are available for computing  $F(\phi, k)$ ,  $E(\phi, k)$ ,  $R_F(a, b, c)$  and  $R_D(a, b, c)$ . DELLPI, DRFVAL, and DRDVAL are double precision routines.

**CALL ELLPI( $\phi, \psi, k, \ell, F, E, IERR$ )**

The arguments  $\phi, \psi, k, \ell$  are real numbers which satisfy  $\phi \geq 0, \psi \geq 0, \phi + \psi = \pi/2$ , and  $k^2 + \ell^2 = 1$ . Also, if  $\psi = 0$  then it is assumed that  $\ell \neq 0$ .  $F, E$ , and  $IERR$  are variables. When ELLPI is called, if no input errors are detected then  $IERR$  is set to 0,  $F$  is assigned the value  $F(\phi, k)$ , and  $E$  is assigned the value  $E(\phi, k)$ .

**Error Return.** If an input error is detected then  $IERR$  is set as follows:

- $IERR = 1$   $\phi < 0$  or  $\psi < 0$
- $IERR = 2$   $|k| > 1$  or  $|\ell| > 1$
- $IERR = 3$   $\psi = 0$  and  $\ell = 0$



**Precision.** ELLPI is accurate to within 4 units of the 14<sup>th</sup> significant digit.

**Programming.** ELLPI calls the functions ALNREL and CPABS. ELLPI was written by Allen V. Hershey and modified by A. H. Morris.

**Reference.** DiDonato, A. R. and Hershey, A. V., "New Formulas for Computing Incomplete Elliptic Integrals of the First and Second Kind," *JACM* 6 (1959), pp. 515-526.

**CALL RFVAL**( $a, b, c, w, \text{IERR}$ )

The arguments  $a, b, c$  are nonnegative real numbers, only one of which can be 0. IERR and  $w$  are variables. When RFVAL is called, if no input errors are detected then IERR is set to 0 and  $w$  is assigned the value  $R_F(a, b, c)$ .

**Error Return.** If an input error is detected then IERR has one of the following values:

IERR = 1 Either  $a, b$ , or  $c$  is negative.

IERR = 2 Either  $a + b, a + c$ , or  $b + c$  is too small.

IERR = 3 Either  $a, b$ , or  $c$  is too large.

**Precision.** RFVAL is accurate to within 4 units of the 14<sup>th</sup> significant digit.

**Programming.** RFVAL was written by B. C. Carlson and Elaine M. Notis (Iowa State University), and modified by A. H. Morris. The function SPMPAR is used.

#### References.

- (1) Carlson, B. C., "Computing Elliptic Integrals by Duplication," *Numerische Mathematik* 33 (1979), pp. 1-16.
- (2) \_\_\_\_\_ and Notis, E. M., "Algorithm 577. Algorithms for Incomplete Elliptic Integrals," *ACM Trans. Math Software* 7 (1981), pp. 398-403.

**CALL RDVAL**( $a, b, c, w, \text{IERR}$ )

The arguments  $a$  and  $b$  are nonnegative real numbers where  $a + b > 0$ , and  $c$  is a positive real number. IERR and  $w$  are variables. When RDVAL is called, if no input errors are detected then IERR is set to 0 and  $w$  is assigned the value  $R_D(a, b, c)$ .

**Error Return.** If an input error is detected then IERR has one of the following values:

IERR = 1 Either  $a, b$ , or  $c$  is negative.

IERR = 2 Either  $a + b$  or  $c$  is too small.

IERR = 3 Either  $a, b$ , or  $c$  is too large.

**Precision.** RDVAL is accurate to within 4 units of the 14<sup>th</sup> significant digit.

**Programming.** RDVAL was written by B. C. Carlson and Elaine M. Notis (Iowa State University), and modified by A. H. Morris. The function SPMPAR is used.

## References.

- (1) Carlson, B. C., "Computing Elliptic Integrals by Duplication," *Numerische Mathematik* 33 (1979), pp. 1-16.
- (2) \_\_\_\_\_ and Notis, E. M., "Algorithm 577, Algorithms for Incomplete Elliptic Integrals," *ACM Trans. Math Software* 7 (1981), pp. 398-403.

### CALL DELLPI( $\phi, \psi, k, \ell, F, E, \text{IERR}$ )

The arguments  $\phi, \psi, k, \ell$  are double precision numbers where  $\phi \geq 0, \psi \geq 0, \phi + \psi = \pi/2$ , and  $k^2 + \ell^2 = 1$ . Also, if  $\psi = 0$  then it is assumed that  $\ell \neq 0$ .  $F$  and  $E$  are double precision variables, and IERR is an integer variable. When DELLPI is called, if no input errors are detected then IERR is set to 0,  $F$  is assigned the double precision value for  $F(\phi, k)$ , and  $E$  is assigned the double precision value for  $E(\phi, k)$ .

**Error Return.** If an input error is detected then IERR is set as follows:

IERR = 1  $\phi < 0$  or  $\psi < 0$   
IERR = 2  $|k| > 1$  or  $|\ell| > 1$   
IERR = 3  $\psi = 0$  and  $\ell = 0$

**Precision.** DELLPI is accurate to within 5 units of the 28<sup>th</sup> significant digit.

**Programming.** DELLPI employs the functions DCPABS, DLNREL, and DPMPAR. DELLPI was written by Allen V. Hershey and modified by A. H. Morris.

**Reference.** DiDonato, A. R. and Hershey, A. V., "New Formulas for Computing Incomplete Elliptic Integrals of the First and Second Kind," *JACM* 6 (1959), pp. 515-526.

### CALL DRFVAL( $a, b, c, w, \text{IERR}$ )

The arguments  $a, b, c$  are nonnegative double precision numbers, only one of which can be 0. IERR is an integer variable and  $w$  a double precision variable. When DRFVAL is called, if no input errors are detected then IERR is set to 0 and  $w$  is assigned the double precision value for  $R_F(a, b, c)$ .

**Error Return.** If an input error is detected then IERR has one of the following values:

IERR = 1 Either  $a, b$ , or  $c$  is negative.  
IERR = 2 Either  $a + b, a + c$ , or  $b + c$  is too small.  
IERR = 3 Either  $a, b$ , or  $c$  is too large.

**Programming.** DRFVAL was written by B. C. Carlson and Elaine M. Notis (Iowa State University), and modified by A. H. Morris. The function DPMPAR is used.

## References.

- (1) Carlson, B. C., "Computing Elliptic Integrals by Duplication," *Numerische Mathematik* 33 (1979), pp. 1-16.

- (2) \_\_\_\_\_ and Notis, E. M., "Algorithm 577, Algorithms for Incomplete Elliptic Integrals," *ACM Trans. Math Software* 7 (1981), pp. 398-403.

**CALL DRDVAL**( $a, b, c, w, \text{IERR}$ )

The arguments  $a$  and  $b$  are nonnegative double precision numbers where  $a + b > 0$ , and  $c$  is a positive double precision number. IERR is an integer variable and  $w$  a double precision variable. When DRDVAL is called, if no input errors are detected then IERR is set to 0 and  $w$  is assigned the double precision value for  $R_D(a, b, c)$ .

**Error Return.** If an input error is detected then IERR has one of the following values:

IERR = 1 Either  $a, b$ , or  $c$  is negative.

IERR = 2 Either  $a + b$  or  $c$  is too small.

IERR = 3 Either  $a, b$ , or  $c$  is too large.

**Programming.** DRDVAL was written by B. C. Carlson and Elaine M. Notis (Iowa State University), and modified by A. H. Morris. The function DPMPAR is used.

**References.**

- (1) Carlson, B. C., "Computing Elliptic Integrals by Duplication," *Numerische Mathematik* 33 (1979), pp. 1-16.
- (2) \_\_\_\_\_ and Notis, E. M., "Algorithm 577, Algorithms for Incomplete Elliptic Integrals," *ACM Trans. Math Software* 7 (1981), pp. 398-403.

## REAL ELLIPTIC INTEGRALS OF THE THIRD KIND

For any  $0 \leq \phi \leq \pi/2$  the elliptic integral  $\Pi(\phi, n, k)$  is defined by

$$\Pi(\phi, n, k) = \int_0^\phi (1 - n \sin^2 \theta)^{-1} (1 - k^2 \sin^2 \theta)^{-\frac{1}{2}} d\theta$$

where  $n$  is any real number such that  $1 - n \sin^2 \phi \neq 0$ , and  $k$  any real number such that  $k^2 \leq 1$  and  $1 - k^2 \sin^2 \phi \neq 0$ . Alternatively, for any  $r \neq 0$  we may consider

$$R_J(a, b, c, r) = \frac{3}{2} \int_0^\infty (t + r)^{-1} [(t + a)(t + b)(t + c)]^{-\frac{1}{2}} dt$$

where  $a, b, c$  are nonnegative and at most one of them is 0. If  $a \leq b \leq c$  and  $a < c$  then

$$R_J(a, b, c, r) = \frac{3c^{-\frac{3}{2}}}{n \sin^3 \phi} [\Pi(\phi, n, k) - F(\phi, k)]$$

where  $F(\phi, k)$  is the elliptic integral of the first kind,  $\cos^2 \phi = a/c$ ,  $k^2 = (c - b)/(c - a)$ , and  $n = (c - r)/(c - a)$ . If  $\phi = \pi/2$  then the elliptic integral  $\Pi(\phi, n, k)$  is said to be *complete*. Otherwise, if  $\phi < \pi/2$  then the integral is said to be *incomplete*. The subroutines EPI, RJVAL, DEPI, and DRJVAL are available for computing  $\Pi(\phi, n, k)$  and  $R_J(a, b, c, r)$ . DEPI and DRJVAL are double precision routines.

**CALL EPI**( $\phi, \psi, k^2, \ell^2, n, m, w, IERR$ )

The arguments  $\phi, \psi, k^2, \ell^2, n, m$  are real numbers where  $\phi \geq 0, \psi \geq 0, \phi + \psi = \pi/2, k^2 + \ell^2 = 1, |n| \leq 1$ , and  $n + m = 1$ . Also, if  $\psi = 0$  then it is assumed that  $\ell^2 \neq 0$  and  $m \neq 0$ . IERR and  $w$  are variables. When EPI is called, if no input errors are detected then IERR is set to 0 and  $w$  is assigned the value  $\Pi(\phi, n, k)$ .

**Error Return.** If an input error is detected then IERR has one of the following values:

IERR = 1 Either  $\phi$  or  $\psi$  is negative, or  $\phi + \psi \neq \pi/2$ .

IERR = 2 Either  $|n| > 1$  or  $n + m \neq 1$ .

IERR = 3 Either  $k^2$  or  $\ell^2$  is negative, or  $k^2 + \ell^2 \neq 1$ .

IERR = 4 Either  $\psi$  and  $m$  are too close to 0, or  $\psi$  and  $\ell^2$  are too close to 0.

**Precision.** EPI is accurate to within 4 units of the 14<sup>th</sup> significant digit.

**Programming.** EPI employs the subroutines RFVAL, RJVAL, RCVAL1 and function SPM-PAR. EPI was written by A.H. Morris.

**CALL RJVAL**( $a, b, c, r, w, IERR$ )

The arguments  $a, b, c$  are nonnegative real numbers, only one of which can be 0, and  $r$  is a positive real number. IERR and  $w$  are variables. When RJVAL is called, if no input

errors are detected then IERR is set to 0 and  $w$  is assigned the value  $R_J(a, b, c, r)$ .

**Error Return.** If an input error is detected then IERR has one of the following values:

- IERR = 1 Either  $a, b, c$ , or  $r$  is negative.
- IERR = 2 Either  $a + b, a + c, b + c$ , or  $r$  is too small.
- IERR = 3 Either  $a, b, c$ , or  $r$  is too large.

**Precision.** RJVAL is accurate to within 4 units of the  $14^{\text{th}}$  significant digit.

**Programming.** RJVAL calls the subroutine RCVAL1. These subroutines were written by B.C. Carlson and Elaine M. Notis (Iowa State University), and modified by A.H. Morris. The function SPMPAR is also used.

#### References.

- (1) Carlson, B. C., "Computing Elliptic Integrals by Duplication." *Numerische Mathematik* 33 (1979), pp. 1-16.
- (2) \_\_\_\_\_ and Notis, E. M., "Algorithm 577, Algorithms for Incomplete Elliptic Integrals." *ACM Trans. Math Software* 7 (1981), pp. 398-403.

**CALL DEPI**( $\phi, \psi, k^2, \ell^2, n, m, w, \text{IERR}$ )

The arguments  $\phi, \psi, k^2, \ell^2, n, m$  are double precision numbers where  $\phi \geq 0, \psi \geq 0, \phi + \psi = \pi/2, k^2 + \ell^2 = 1, |n| \leq 1$ , and  $n + m = 1$ . Also, if  $\psi = 0$  then it is assumed that  $\ell^2 \neq 0$  and  $m \neq 0$ . IERR is an integer variable and  $w$  a double precision variable. When DEPI is called, if no input errors are detected then IERR is set to 0 and  $w$  is assigned the double precision value for  $\Pi(\phi, n, k)$ .

**Error Return.** If an input error is detected then IERR has one of the following values:

- IERR = 1 Either  $\phi$  or  $\psi$  is negative, or  $\phi + \psi \neq \pi/2$ .
- IERR = 2 Either  $|n| > 1$  or  $n + m \neq 1$ .
- IERR = 3 Either  $k^2$  or  $\ell^2$  is negative, or  $k^2 + \ell^2 \neq 1$ .
- IERR = 4 Either  $\psi$  and  $m$  are too close to 0, or  $\psi$  and  $\ell^2$  are too close to 0.

**Programming.** DEPI employs the subroutines DRFVAL, DRJVAL, DRCVL1 and function DPMPAR. DEPI was written by A.H. Morris.

**CALL DRJVAL**( $a, b, c, r, w, \text{IERR}$ )

The arguments  $a, b, c$  are nonnegative double precision numbers, only one of which can be 0, and  $r$  is a positive double precision number. IERR is an integer variable and  $w$  a double precision variable. When DRJVAL is called, if no input errors are detected then IERR is set to 0 and  $w$  is assigned the double precision value for  $R_J(a, b, c, r)$ .

**Error Return.** If an input error is detected then IERR has one of the following values:

- IERR = 1 Either  $a, b, c$ , or  $r$  is negative.
- IERR = 2 Either  $a + b, a + c, b + c$ , or  $r$  is too small.

IERR = 3 Either  $a, b, c$ , or  $r$  is too large.

**Programming.** DRJVAL calls the subroutine DRCVL1. These subroutines were written by B.C. Carlson and Elaine M. Notis (Iowa State University), and modified by A.H. Morris. The function DPMPAR is also used.

**References.**

- (1) Carlson, B. C., "Computing Elliptic Integrals by Duplication." *Numerische Mathematik* **33** (1979), pp. 1-16.
- (2) \_\_\_\_\_ and Notis, E. M., "Algorithm 577, Algorithms for Incomplete Elliptic Integrals." *ACM Trans. Math Software* **7** (1981), pp.398-403.

## JACOBIAN ELLIPTIC FUNCTIONS

For any complex number  $k \neq 0, \pm 1$  the elliptic function  $sn(z, k)$  may be defined as the meromorphic function  $w(z)$  that satisfies

$$(1) \quad \left( \frac{dw}{dz} \right)^2 = (1 - w^2)(1 - k^2 w^2)$$

$$w(0) = 0, \quad w'(0) = 1.$$

If  $k = 0$  then  $sn(z, 0) = \sin z$  satisfies (1), and if  $k = \pm 1$  then  $sn(z, k) = \tanh z$  satisfies (1). Alternatively,  $sn(z, k) = \sin \phi$  where  $\phi(z)$  satisfies

$$(2) \quad \left( \frac{d\phi}{dz} \right)^2 = 1 - k^2 \sin^2 \phi$$

$$\phi(0) = 0, \quad \phi'(0) = 1.$$

The elliptic functions  $cn(z, k)$  and  $dn(z, k)$  may be defined by

$$cn(z, k) = \sqrt{1 - sn^2(z, k)}$$

$$dn(z, k) = \sqrt{1 - k^2 sn^2(z, k)}$$

where the roots take the value 1 for  $z = 0$ . In particular, if  $k = 0$  then  $cn(z, 0) = \cos z$  and  $dn(z, 0) = 1$ , and if  $k = \pm 1$  then  $cn(z, k) = dn(z, k) = 1/\cosh z$ . The subroutines ELLPF and ELPFC1 are available for computing  $sn(z, k)$ ,  $cn(z, k)$ , and  $dn(z, k)$  when  $k$  is a real value such that  $|k| \leq 1$ . ELLPF may be used when  $z$  is real and ELPFC1 when  $z$  is complex.

### CALL ELLPF( $u, k, \ell, S, C, D, IERR$ )

It is assumed that  $u, k$ , and  $\ell$  are real numbers where  $k^2 + \ell^2 = 1$ .  $S, C$ , and  $D$  are real variables. When ELLPF is called,  $S, C$ , and  $D$  are assigned the values  $S = sn(u, k)$ ,  $C = cn(u, k)$ , and  $D = dn(u, k)$ .

IERR is a variable that reports the status of the results. When the routine terminates, IERR has one of the following values:

IERR = 0 The elliptic functions were computed.

IERR = 1 (Input error)  $k^2 + \ell^2 \neq 1$ .

IERR = 2  $u$  is too large for  $k$ .

When  $IERR \geq 1$ , no computation is performed.

**Precision.** Let  $K(k)$  be the complete elliptic integral of the first kind. For  $|k| < .99995$  the relative errors of  $sn(u, k)$  and  $dn(u, k)$  are less than  $10^{-12}$  when  $0 < u \leq K(k)$ , and the relative error of  $cn(u, k)$  is less than  $10^{-12}$  when  $0 \leq u \leq .97K(k)$ .

**Algorithm.** Let  $K = K(k)$  be the complete elliptic integral of the first kind. For  $0 \leq u \leq K/2$  (when  $\ell \neq 0$ ), the Maclaurin expansion

$$sn(u, k) = u - \frac{(1+k^2)u^3}{3!} + \frac{(1+14k^2+k^4)u^5}{5!} - \dots$$

is employed when  $u \leq .01$ . Otherwise, if  $u > .01$  let  $K' = K(\ell)$ ,  $q = \exp(-\pi K'/K)$ , and  $r = \exp(-\pi K/K')$ . Then

$$sn(u, k) = \frac{2\pi}{kK} \sum_{n \geq 0} \frac{q^{n+\frac{1}{2}}}{1-q^{2n+1}} \sin \frac{(2n+1)\pi u}{2K}$$

is used when  $k \leq \ell$  and

$$sn(u, k) = \frac{\pi}{2kK'} \left[ \tanh \frac{\pi u}{2K'} + 4 \sum_{n \geq 1} \frac{(-1)^n r^{2n}}{1+r^{2n}} \sinh \frac{2\pi u}{K'} \right]$$

is used when  $k > \ell$ . The functions  $cn(u, k)$  and  $dn(u, k)$  are obtained from

$$\begin{aligned} sn(u, k)^2 + cn(u, k)^2 &= 1 \\ dn(u, k)^2 + k^2 sn(u, k)^2 &= 1. \end{aligned}$$

For  $K/2 < u \leq K$  the identities

$$\begin{aligned} sn(u, k) &= cn(v, k)/dn(v, k) \\ cn(u, k) &= |\ell| sn(v, k)/dn(v, k) \\ dn(u, k) &= |\ell|/dn(v, k) \end{aligned}$$

are applied. Here  $v = K - u$ .

**Programming.** ELLPF employs the subroutines SCD, SCDF, SCDJ, SCDM, ELLPI, SNHCSH and functions ALNREL, CPABS, SPMPAR, IPMPAR. ELLPF was written by Andrew H. van Tuyl and modified by A. H. Morris.

**CALL ELPFC1**( $z, k, \ell, S, C, D, IERR$ )

The argument  $z$  is complex, and  $k$  and  $\ell$  are real numbers where  $k^2 + \ell^2 = 1$ .  $S, C$ , and  $D$  are complex variables. When ELPFC1 is called  $S, C$ , and  $D$  are assigned the values  $S = sn(z, k)$ ,  $C = cn(z, k)$ , and  $D = dn(z, k)$ .

IERR is a variable that reports the status of the results. When the routine terminates, IERR has one of the following values:

- IERR = 0 The elliptic functions were computed.
- IERR = 1 (Input error)  $k^2 + \ell^2 \neq 1$ .
- IERR = 2  $z$  is too large for  $k$ .
- IERR = 3  $z$  is a pole for the elliptic functions.



When  $IERR \geq 1$ , no computation is performed.

**Precision.** Let  $x = \text{Re}(z)$ ,  $y = \text{Im}(z)$ ,  $K = K(k)$  be the elliptic integral of the first kind for  $k$ , and  $K' = K(\ell)$ . For  $|k| < .99995$  the relative errors of the real and imaginary parts of  $sn(z, k)$  are less than  $10^{-12}$  when  $0 < x \leq K$  and  $0 < y \leq .992K'$ , and the relative errors of the real and imaginary parts of  $cn(z, k)$  and  $dn(z, k)$  are less than  $10^{-12}$  when  $0 < x \leq .97K$  and  $0 < y \leq .97K'$ .

**Algorithm.** For  $z = x + iy$  let

$$\begin{array}{ll} s = sn(x, k) & s_1 = sn(y, \ell) \\ c = cn(x, k) & c_1 = cn(y, \ell) \\ d = dn(x, k) & d_1 = dn(y, \ell) \end{array}$$

and  $D = c_1^2 + k^2 s^2 s_1^2$ . Then

$$\begin{array}{l} sn(z, k) = (sd_1 + icds_1c_1)/D \\ cn(z, k) = (cc_1 - isds_1d_1)/D \\ dn(z, k) = (dc_1d_1 - ik^2scs_1)/D \end{array}$$

are applied when  $D \neq 0$ .

**Programming.** ELPFC1 calls ELLPF, which employs the subroutines SCD, SCDF, SCDJ, SCDM, ELLPI, SNHCSH and functions ALNREL, CPABS, SPMPAR, IPMPAR. ELPFC1 was written by Andrew H. van Tuyl.

## WEIERSTRASS ELLIPTIC FUNCTION FOR THE EQUIANHARMONIC AND LEMNISCATIC CASES

Let  $w$  and  $w'$  be complex numbers where  $\text{Im}(w'/w) > 0$ , and  $w_{mn} = 2mw + 2nw'$  for all integers  $m, n$ . Then for any complex  $z$ , the Weierstrass elliptic function  $\mathcal{P}(z; w, w')$  can be defined by

$$\mathcal{P}(z; w, w') = \frac{1}{z^2} + \sum' \left[ \frac{1}{(z - w_{mn})^2} - \frac{1}{w_{mn}^2} \right]$$

where  $\Sigma'$  denotes the sum for all  $m, n = 0, \pm 1, \pm 2, \dots$  except  $m = n = 0$ . If  $w = re^{i\phi}$  and  $w' = r'e^{i\phi'}$  where  $\phi' = \phi + \theta$  for  $0 \leq \theta < 2\pi$ , then the restriction  $\text{Im}(w'/w) > 0$  is equivalent to assuming that  $0 < \theta < \pi$ .  $\mathcal{P}(z; w, w')$  is analytic everywhere except at the points  $w_{mn}$ , which are poles, and

$$\mathcal{P}(z + 2w; w, w') = \mathcal{P}(z; w, w')$$

$$\mathcal{P}(z + 2w'; w, w') = \mathcal{P}(z; w, w')$$

for all  $z$ . The relations

$$\mathcal{P}(-z; w, w') = \mathcal{P}(z; w, w')$$

$$\mathcal{P}(\lambda z; \lambda w, \lambda w') = \lambda^{-2} \mathcal{P}(z; w, w') \quad \lambda \neq 0$$

also hold. A somewhat surprising fact is that only the values  $g_2 = 60\Sigma' w_{mn}^{-4}$  and  $g_3 = 140\Sigma' w_{mn}^{-6}$  are needed for computing  $\mathcal{P}(z; w, w')$  at a point  $z$ . Hence,  $\mathcal{P}(z; w, w')$  is frequently denoted by  $\mathcal{P}(z; g_2, g_3)$ . For  $\lambda \neq 0$

$$g_2(\lambda w, \lambda w') = \lambda^{-4} g_2(w, w')$$

$$g_3(\lambda w, \lambda w') = \lambda^{-6} g_3(w, w')$$

also hold. We now consider the following cases:

(1) Equianharmonic ( $g_2 = 0$  and  $g_3$  is a positive real number)

(2) Lemniscatic ( $g_2$  is a positive real number and  $g_3 = 0$ )

(1) occurs when  $2w = \frac{1}{2} - \frac{\sqrt{3}}{2}i$  and  $2w' = \frac{1}{2} + \frac{\sqrt{3}}{2}i$ , and (2) occurs when  $2w = 1$  and  $2w' = i$ . The following subroutines are available for computing  $\mathcal{P}(z; w, w')$  and its derivative  $\mathcal{P}'(z; w, w')$  for these two choices of  $(w, w')$ .

**CALL PEQ( $z, e, \text{IERR}$ )**

The argument  $z$  is a complex number,  $e$  is a complex variable, and **IERR** is an integer variable. It is assumed that the periods are  $2w = \frac{1}{2} - \frac{\sqrt{3}}{2}i$  and  $2w' = \frac{1}{2} + \frac{\sqrt{3}}{2}i$ . When **PEQ** is called, if  $z$  is not a pole then **IERR** is assigned the value 0 and  $e$  is assigned the value  $\mathcal{P}(z; w, w')$ .

**Error Return.** If  $z = w_{mn}$  for some  $m, n$  then **IERR** is assigned the value 1 and  $e = 0$ .

**Precision.** If  $|\mathcal{P}(z; w, w')| \leq 1$  then the absolute error is less than  $7 \cdot 10^{-13}$ . Otherwise, the relative error is less than  $7 \cdot 10^{-13}$ .

**Programming.** Written by Ulrich Eckhardt (University of Hamburg, West Germany). Modified by A. H. Morris.

#### References.

- (1) Eckhardt, Ulrich, "Algorithm 549, Weierstrass' Elliptic Functions," *ACM Trans. Math Software* 4 (1980), pp.112-120.
- (2) \_\_\_\_\_, "A Rational Approximation to Weierstrass'  $\mathcal{P}$ -Function," *Math Comp.* 30 (1976), pp.818-826.

#### CALL PEQ1( $z, e, IERR$ )

The argument  $z$  is a complex number,  $e$  is a complex variable, and  $IERR$  is an integer variable. It is assumed that the periods are  $2w = \frac{1}{2} - \frac{\sqrt{3}}{2}i$  and  $2w' = \frac{1}{2} + \frac{\sqrt{3}}{2}i$ . When PEQ1 is called, if  $z$  is not a pole then  $IERR$  is assigned the value 0 and  $e$  is assigned the value  $\mathcal{P}'(z; w, w')$ .

**Error Return.** If  $z = w_{mn}$  for some  $m, n$  then  $IERR$  is assigned the value 1 and  $e = 0$ .

**Precision.** If  $|\mathcal{P}'(z; w, w')| \leq 1$  then the absolute error is less than  $7 \cdot 10^{-13}$ . Otherwise, the relative error is less than  $7 \cdot 10^{-13}$ .

**Programming.** Written by Ulrich Eckhardt (University of Hamburg, West Germany). Modified by A. H. Morris.

#### References.

- (1) Eckhardt, Ulrich, "Algorithm 549, Weierstrass' Elliptic Functions," *ACM Trans. Math Software* 4 (1980), pp.112-120.
- (2) \_\_\_\_\_, "A Rational Approximation to Weierstrass'  $\mathcal{P}$ -Function," *Math Comp.* 30 (1976), pp.818-826.

#### CALL PLEM( $z, e, IERR$ )

The argument  $z$  is a complex number,  $e$  is a complex variable, and  $IERR$  is an integer variable. It is assumed that the periods are  $2w = 1$  and  $2w' = i$ . When PLEM is called, if  $z$  is not a pole then  $IERR$  is assigned the value 0 and  $e$  is assigned the value  $\mathcal{P}(z; w, w')$ .

**Error Return.** If  $z = w_{mn}$  for some  $m, n$  then  $IERR$  is assigned the value 1 and  $e = 0$ .

**Precision.** If  $|\mathcal{P}(z; w, w')| \leq 1$  then the absolute error is less than  $6 \cdot 10^{-13}$ . Otherwise, the relative error is less than  $6 \cdot 10^{-13}$ .

**Programming.** Written by Ulrich Eckhardt (University of Hamburg, West Germany). Modified by A. H. Morris.

#### References.

- (1) Eckhardt, Ulrich, "Algorithm 549, Weierstrass' Elliptic Functions," *ACM Trans. Math Software* 4 (1980), pp.112-120.

- (2) \_\_\_\_\_, "A Rational Approximation to Weierstrass'  $\mathcal{P}$ -Function. II: The Lemniscatic Case," *Computing (Arch. Elektron. Rechnen)* 18 (1977), pp. 341-349.

**CALL PLEM1( $z, e, IERR$ )**

The argument  $z$  is a complex number,  $e$  is a complex variable, and  $IERR$  is an integer variable. It is assumed that the periods are  $2w = 1$  and  $2w' = i$ . When PLEM1 is called, if  $z$  is not a pole then  $IERR$  is assigned the value 0 and  $e$  is assigned the value  $\mathcal{P}'(z; w, w')$ .

**Error Return.** If  $z = w_{mn}$  for some  $m, n$  then  $IERR$  is assigned the value 1 and  $e = 0$ .

**Precision.** If  $|\mathcal{P}'(z; w, w')| \leq 1$  then the absolute error is less than  $6 \cdot 10^{-13}$ . Otherwise, the relative error is less than  $6 \cdot 10^{-13}$ .

**Programming.** Written by Ulrich Eckhardt (University of Hamburg, West Germany). Modified by A. H. Morris.

**References.**

- (1) Eckhardt, Ulrich, "Algorithm 549, Weierstrass' Elliptic Functions," *ACM Trans. Math Software* 4 (1980), pp.112-120.
- (2) \_\_\_\_\_, "A Rational Approximation to Weierstrass'  $\mathcal{P}$ -Function. II: The Lemniscatic Case," *Computing (Arch. Elektron. Rechnen)* 18 (1977), pp. 341-349.

## INTEGRAL OF THE BIVARIATE DENSITY FUNCTION OVER ARBITRARY POLYGONS AND SEMI-INFINITE ANGULAR REGIONS

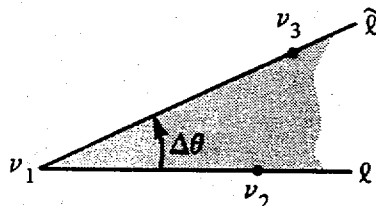
Given a sequence of points  $\nu_i = (x_i, y_i)$  ( $i = 1, \dots, n+1$ ) where  $n \geq 3$  and  $\nu_{n+1} = \nu_1$ . Let  $\tau$  denote the polygon whose boundary  $\partial\tau$  is a polygonal line which begins at point  $\nu_1$ , traverses the points  $\nu_i$  in the order that they are indexed, and is the straight line segment connecting  $\nu_i$  to  $\nu_{i+1}$  for each  $i = 1, \dots, n$  where  $\nu_i \neq \nu_{i+1}$ . Then the subroutine VALR2 is available for computing the integral

$$P(\tau) = \frac{1}{2\pi} \iint_{\tau} e^{-(x^2+y^2)/2} dx dy$$

and the associated function  $A(\tau) = \iint_{\tau} dx dy$ . If the boundary  $\partial\tau$  is a simple positively (negatively) oriented closed curve, then  $P(\tau)$  and  $A(\tau)$  are positive (negative) and  $|A(\tau)| =$  the area of  $\tau$ . However,  $\partial\tau$  need not be simple. It may be self-intersecting or have overlapping line segments. If  $\Delta\theta_i$  is the angle between the vectors  $\nu_i - \nu_{i-1}$  and  $\nu_{i+1} - \nu_i$  (where  $\nu_0 = \nu_n$ ), then it may occur that  $\Delta\theta_i = \pi$  for some  $i$ , in which case a portion of the polygon may be degenerate. In general,  $-\pi < \Delta\theta_i \leq \pi$  for each  $i$  where the sign of the angle is positive (negative) if the angle is measured in a counterclockwise (clockwise) direction from  $\nu_i - \nu_{i-1}$  to  $\nu_{i+1} - \nu_i$ . VALR2 also computes the value  $k(\tau) = \frac{1}{2\pi} \sum_{i=1}^n \Delta\theta_i$ , which is an integer. If the boundary is a simple closed curve, then  $k(\tau)$  is the winding number of the curve around any interior point of the polygon  $\tau$ .

Alternatively, assume that we are given three points  $\nu_i = (x_i, y_i)$  ( $i = 1, 2, 3$ ) and let  $\Delta\theta$  denote the angle between the vectors  $\nu_2 - \nu_1$  and  $\nu_3 - \nu_1$ . In this case, assume that the angle  $\Delta\theta$  is measured in a counterclockwise direction from  $\nu_2 - \nu_1$  to  $\nu_3 - \nu_1$ , so that  $0 \leq \Delta\theta < 2\pi$ . Let  $\ell$  denote the straight line beginning at point  $\nu_1$  and passing through point  $\nu_2$ , and let  $\tilde{\ell}$  denote the straight line beginning at  $\nu_1$  and passing through  $\nu_3$ . Then the subroutine VALR2 is also available for computing  $P(\tau)$  when  $\tau$  is the semi-infinite angular region bounded by  $\ell$  and  $\tilde{\ell}$ , and having the angle  $\Delta\theta$ .  $0 \leq P(\tau) \leq 1$  for any angular region  $\tau$ , and  $P(\tau) \rightarrow 1$  when  $\Delta\theta \rightarrow 2\pi$ .

Angular region  $\tau$



CALL VALR2( $X, Y, n, P, IOP, A, IND, k$ )

The argument  $n$  is either 1 or the number of points involved in defining a polygon. If  $n = 1$  then it is assumed that  $\tau$  is a semi-infinite angular region defined by the points  $\nu_i = (x_i, y_i)$  ( $i = 1, 2, 3$ ), and that  $X$  and  $Y$  are arrays containing  $x_1, x_2, x_3$  and  $y_1, y_2, y_3$ . Otherwise, if  $n \neq 1$  then it is assumed that  $\tau$  is a polygon defined by the points  $\nu_i =$

$(x_i, y_i)$  ( $i = 1, \dots, n+1$ ) where  $n \geq 3$  and  $\nu_{n+1} = \nu_1$ . In this case,  $X$  and  $Y$  are arrays containing the abscissae  $x_1, \dots, x_{n+1}$  and ordinates  $y_1, \dots, y_{n+1}$ . Since  $\nu_{n+1} = \nu_1$ , the values  $x_{n+1}$  and  $y_{n+1}$  need not be supplied by the user. The routine automatically stores  $x_1$  and  $y_1$  in  $X(n+1)$  and  $Y(n+1)$ .

$P$ ,  $A$ , and  $k$  are variables. If  $n = 1$  then  $P$  is assigned the value  $P(\tau)$  for the angular region  $\tau$  and  $A$  is assigned the value 0. In this case,  $k$  is not defined. Otherwise, if  $n \geq 3$  then  $P$  is assigned the value  $P(\tau)$ ,  $A$  is assigned the value  $A(\tau)$ , and  $k$  is assigned the value  $k(\tau)$  for the polygon  $\tau$ .

IOP is an input argument which specifies the (relative) precision to which  $P(\tau)$  is to be computed. IOP is set to 1, 2, or 3 for 3, 6, or 9 decimal digit accuracy.

IND is the variable that reports the status of the results. The routine assigns IND one of the following values:

- IND = 0 The desired values were obtained.
- IND = 1 (Input error)  $\nu_1$  is either equal to  $\nu_2$  or  $\nu_3$ , or is too close to  $\nu_2$  or  $\nu_3$  to compute  $P(\tau)$  for the angular region  $\tau$ . In this case,  $P$  is set to 5.
- IND = 2 The desired values were obtained. If  $n = 1$  then  $\Delta\theta \approx \pi$ . Otherwise, if  $n \geq 3$  then  $|\Delta\theta_i| \approx \pi$  for some  $i$ .
- IND = 3 (Input error) Either  $n < 1$  or  $n = 2$ .

**Remarks.** VALR2 can be used for computing the integral of the general bivariate density function over an arbitrary polygon or semi-infinite angular region  $\hat{\tau}$ . Consider

$$\hat{P}(\hat{\tau}) = B \iint_{\hat{\tau}} \exp \left\{ \frac{-1}{2(1-\rho^2)} \left[ \left( \frac{\omega - \mu_\omega}{\sigma_\omega} \right)^2 - 2\rho \frac{(\omega - \mu_\omega)(z - \mu_z)}{\sigma_\omega \sigma_z} + \left( \frac{z - \mu_z}{\sigma_z} \right)^2 \right] \right\} d\omega dz$$

where  $B = (1 - \rho^2)^{-1/2} / (2\pi\sigma_\omega\sigma_z)$ ,  $(\mu_\omega, \mu_z)$  is the mean,  $\sigma_\omega$  and  $\sigma_z$  are the (nonzero) variances, and  $\rho$  is the correlation coefficient satisfying  $|\rho| < 1$ . Consider also

$$x = (1 - \rho^2)^{-\frac{1}{2}} \left[ \frac{\omega - \mu_\omega}{\sigma_\omega} - \rho \frac{z - \mu_z}{\sigma_z} \right] \quad \text{and} \quad y = \frac{z - \mu_z}{\sigma_z}.$$

Since this transformation maps straight lines into straight lines,  $\hat{\tau}$  is mapped onto a polygon or angular region  $\tau$  and we obtain  $\hat{P}(\hat{\tau}) = P(\tau)$ . Moreover, if  $\hat{\tau}$  is a polygon then  $A(\hat{\tau}) = \iint_{\hat{\tau}} d\omega dz = \sigma_\omega \sigma_z \sqrt{1 - \rho^2} A(\tau)$ .

**Programming.** VALR2 employs the functions ERF, ERFC1, and SPMPAR. VALR2 was designed by Armido R. DiDonato and Richard K. Hageman, and modified by A. H. Morris.

**Reference.** DiDonato, A. R., and Hageman, R. K., *Computation of the Integral of the Bivariate Normal Distribution over Arbitrary Polygons*, Report TR 80-166, Naval Surface Weapons Center, Dahlgren, Virginia, 1980.

## CIRCULAR COVERAGE FUNCTION

The subroutine CIRC is available for computing the circular coverage function  $P(R, d)$  and the generalized circular error function  $V(K, c)$ .  $V$  is the integral of an uncorrelated elliptical Gaussian distribution with standard deviations  $\sigma_x$  and  $\sigma_y$  over a circle of radius  $K\sigma_x$  centered at the mean of the distribution. If  $\sigma_x \geq \sigma_y$  then

$$V(K, c) = \frac{1}{\pi c} \int_0^K \int_0^\pi \exp \left\{ \frac{-r^2}{4c^2} [1 + c^2 + (1 - c^2) \cos \theta] \right\} r dr d\theta$$

where  $c = \sigma_y/\sigma_x$ .  $P$  is the integral of a circular Gaussian distribution with common standard deviation  $\sigma$  over a circle of radius  $R\sigma$  whose center is offset a distance  $d\sigma$  from the mean of the distribution.

$$P(R, d) = \frac{1}{2\pi} \int_0^R \int_0^{2\pi} \exp \left\{ -\frac{1}{2} [(d + r \cos \theta)^2 + r^2 \sin^2 \theta] \right\} r dr d\theta$$

**CALL CIRC( $x, a, i, w, IERR$ )**

The argument  $i$  may be any integer. If  $i = 0$  then the arguments  $x$  and  $a$  are assumed to have the values  $x = K$  and  $a = c$  where  $K \geq 0$  and  $0 \leq c \leq 1$ . Otherwise, if  $i \neq 0$  then  $x = R$  and  $a = d$  where  $R \geq 0$  and  $d \geq 0$ .

$IERR$  and  $w$  are variables. When CIRC is called, if no input errors are detected then  $IERR$  is assigned the value 0. Also,  $w = V(K, c)$  if  $i = 0$  and  $w = P(R, d)$  if  $i \neq 0$ .

**Error Return.** If an input error is detected then  $IERR$  is set as follows:

$IERR = 1$   $x \geq 0$  is not satisfied.

$IERR = 2$   $0 \leq c \leq 1$  or  $d \geq 0$  is not satisfied.

When either of these errors is detected,  $w$  is assigned the value  $-1$ .

**Precision.** CIRC is accurate to within  $10^{-6}$ .

**Note.** If  $\sigma_x \leq \sigma_y$  then reverse the roles of  $x$  and  $y$ .

**Programming.** CIRC calls these functions ERF0 and ERFC0. The routine is an adaptation by A. H. Morris of the BASIC program CIRC given in reference (1).

### References.

- (1) DiDonato, A. R., *Five Statistical Programs in BASIC for Desktop Computers*, Report NSWC TR 83-13, Naval Surface Weapons Center, Dahlgren, Virginia, 1982.
- (2) \_\_\_\_\_ and Jarnagin, M. P., *A Method for Computing the Generalized Circular Error Function and the Circular Coverage Function*, Report 1768, Naval Weapons Laboratory, Dahlgren, Virginia, 1962.

## ELLIPTICAL COVERAGE FUNCTION

The subroutines PKILL and PKILL3 are available for evaluating the integral of an uncorrelated elliptical Gaussian distribution over the area  $A$  of a circle  $(x-h)^2 + (y-k)^2 = R^2$ . The probability to be computed is given by

$$P(R, \sigma_x, \sigma_y, h, k) = \frac{1}{2\pi\sigma_x\sigma_y} \iint_A \exp \left[ -\frac{1}{2} \left( \frac{x^2}{\sigma_x^2} + \frac{y^2}{\sigma_y^2} \right) \right] dx dy$$

where  $\sigma_x$  is the standard deviation in the  $x$  direction and  $\sigma_y$  is the standard deviation in the  $y$  direction.

CALL PKILL( $R, \sigma_x, \sigma_y, h, k, p$ )  
CALL PKILL3( $R, \sigma_x, \sigma_y, h, k, p$ )

$R, \sigma_x, \sigma_y, h, k$  are real numbers and  $p$  is a variable. It is assumed that  $R \geq 0$ ,  $\sigma_x > 0$ , and  $\sigma_y > 0$ . When PKILL or PKILL3 is called,  $p$  is assigned the value  $P(R, \sigma_x, \sigma_y, h, k)$ .

**Precision.** PKILL is accurate to within  $10^{-6}$  and PKILL3 is accurate to within  $10^{-3}$ .

**Programming.** PKILL and PKILL3 call the function ERFC2. The routines are adaptations by A. H. Morris of the BASIC programs ELLCV and ELLCV3 given in reference (1).

### References.

- (1) DiDonato, A. R., *Five Statistical Programs in BASIC for Desktop Computers*, Report TR 83-13, Naval Surface Weapons Center, Dahlgren, Virginia, 1982.
- (2) \_\_\_\_\_ and Jarnagin, M. P., *Integration of the General Bivariate Distribution Over an Offset Ellipse*, Report 1710, Naval Weapons Laboratory, Dahlgren, Virginia, 1960.



## COPYING POLYNOMIALS

If  $p(x) = \sum_{j=0}^{m-1} a_j x^j$  and the coefficients  $a_j$  are stored in an array  $A$ , then the following subroutines are available for copying the first  $n$  coefficients  $a_j$  into an array  $B$ .

**CALL PLCOPY( $A, ka, m, B, kb, n$ )**  
**CALL DPCOPY( $A, ka, m, B, kb, n$ )**

$A$  and  $B$  are arrays. PLCOPY is used if  $A$  and  $B$  are real arrays, and DPCOPY is used if  $A$  and  $B$  are double precision arrays.

The arguments  $m, n, ka, kb$  are positive integers. The coefficients  $a_j$  are assumed to be stored in  $A$  where  $A(1 + j \cdot ka) = a_j$  for  $j = 0, 1, \dots, m - 1$ . The routine stores the first  $n$  coefficients  $a_j$  in  $B$  where  $B(1 + j \cdot kb) = a_j$  for  $j = 0, 1, \dots, n - 1$ .

**Note.** If  $n > m$  then  $B(1 + j \cdot kb) = 0$  for  $j \geq m$ .

**Programmer.** A. H. Morris

## ADDITION OF POLYNOMIALS

If  $p(x) = \sum_{j=0}^{\ell-1} a_j x^j$  and  $q(x) = \sum_{j=0}^{m-1} b_j x^j$  then the following subroutines are available for computing the first  $n$  coefficients of the polynomial  $p(x) + q(x) = \sum_j c_j x^j$ .

**CALL PADD**( $A, ka, \ell, B, kb, m, C, kc, n$ )  
**CALL DPADD**( $A, ka, \ell, B, kb, m, C, kc, n$ )

$A$ ,  $B$  and  $C$  are arrays. PADD is used if  $A$ ,  $B$  and  $C$  are real arrays and DPADD is used if  $A$ ,  $B$  and  $C$  are double precision arrays.

The arguments  $\ell, m, n, ka, kb, kc$  are positive integers. The coefficients  $a_j$  and  $b_j$  are assumed to be stored in  $A$  and  $B$  where

$$A(1 + j \cdot ka) = a_j \quad (j = 0, 1, \dots, \ell - 1)$$

$$B(1 + j \cdot kb) = b_j \quad (j = 0, 1, \dots, m - 1).$$

The routine stores the first  $n$  coefficients  $c_j$  of  $p(x) + q(x)$  in  $C$  where  $C(1 + j \cdot kc) = c_j$  for  $j = 0, 1, \dots, n - 1$ .

**Remarks.** The array  $C$  may begin in the same location as  $A$  or  $B$ . If  $C$  begins in the same location as  $A$  then it is assumed that  $kc = ka$ . In this case, the result  $C$  will overwrite the input data  $A$ . Similarly, if  $C$  begins in the same location as  $B$  then it is assumed that  $kc = kb$ . Otherwise, if  $C$  does not begin in the same location as  $A$  or  $B$ , then it is assumed that the array  $C$  does not overlap with the arrays  $A$  and  $B$ .

**Programmer.** A. H. Morris

## SUBTRACTION OF POLYNOMIALS

If  $p(x) = \sum_{j=0}^{\ell-1} a_j x^j$  and  $q(x) = \sum_{j=0}^{m-1} b_j x^j$  then the following subroutines are available for computing the first  $n$  coefficients of the polynomial  $p(x) - q(x) = \sum_j c_j x^j$ .

**CALL PSUBT**( $A, ka, \ell, B, kb, m, C, kc, n$ )  
**CALL DPSUBT**( $A, ka, \ell, B, kb, m, C, kc, n$ )

$A, B$  and  $C$  are arrays. PSUBT is used if  $A, B$  and  $C$  are real arrays and DPSUBT is used if  $A, B$  and  $C$  are double precision arrays.

The arguments  $\ell, m, n, ka, kb, kc$  are positive integers. The coefficients  $a_j$  and  $b_j$  are assumed to be stored in  $A$  and  $B$  where

$$\begin{aligned} A(1 + j \cdot ka) &= a_j \quad (j = 0, 1, \dots, \ell - 1) \\ B(1 + j \cdot kb) &= b_j \quad (j = 0, 1, \dots, m - 1). \end{aligned}$$

The routine stores the first  $n$  coefficients  $c_j$  of  $p(x) - q(x)$  in  $C$  where  $C(1 + j \cdot kc) = c_j$  for  $j = 0, 1, \dots, n - 1$ .

**Remarks.** The array  $C$  may begin in the same location as  $A$  or  $B$ . If  $C$  begins in the same location as  $A$  then it is assumed that  $kc = ka$ . In this case, the result  $C$  will overwrite the input data  $A$ . Similarly, if  $C$  begins in the same location as  $B$  then it is assumed that  $kc = kb$ . Otherwise, if  $C$  does not begin in the same location as  $A$  or  $B$ , then it is assumed that the array  $C$  does not overlap with the arrays  $A$  and  $B$ .

**Programmer.** A. H. Morris

## MULTIPLICATION OF POLYNOMIALS

If  $p(x) = \sum_{j=0}^{\ell-1} a_j x^j$  and  $q(x) = \sum_{j=0}^{m-1} b_j x^j$  then the following subroutines are available for computing the first  $n$  coefficients of the polynomial  $p(x)q(x) = \sum_j c_j x^j$ .

**CALL PMULT**( $A, ka, \ell, B, kb, m, C, kc, n$ )  
**CALL DPMULT**( $A, ka, \ell, B, kb, m, C, kc, n$ )

$A, B$  and  $C$  are arrays. PMULT is used if  $A, B$  and  $C$  are real arrays and DPMULT is used if  $A, B$  and  $C$  are double precision arrays.

The arguments  $\ell, m, n, ka, kb, kc$  are positive integers. The coefficients  $a_j$  and  $b_j$  are assumed to be stored in  $A$  and  $B$  where

$$A(1 + j \cdot ka) = a_j \quad (j = 0, 1, \dots, \ell - 1)$$

$$B(1 + j \cdot kb) = b_j \quad (j = 0, 1, \dots, m - 1).$$

The routine stores the first  $n$  coefficients  $c_j$  of  $p(x)q(x)$  in  $C$  where  $C(1 + j \cdot kc) = c_j$  for  $j = 0, 1, \dots, n - 1$ .

**Remarks.** It is assumed that the array  $C$  does not overlap with the arrays  $A$  and  $B$ .

**Programmer.** A. H. Morris

## DIVISION OF POLYNOMIALS

If  $p(x) = \sum_{j=0}^{\ell-1} a_j x^j$  and  $q(x) = \sum_{j=0}^{m-1} b_j x^j$  where  $b_0 \neq 0$ , then the following subroutines are available for computing the first  $n$  coefficients of the series  $p(x)/q(x) = \sum_j c_j x^j$ .

**CALL PDIV( $A, ka, \ell, B, kb, m, C, kc, n, IERR$ )**  
**CALL DPDIV( $A, ka, \ell, B, kb, m, C, kc, n, IERR$ )**

$A$ ,  $B$  and  $C$  are arrays. PDIV is used if  $A$ ,  $B$  and  $C$  are real arrays and DPDIV is used if  $A$ ,  $B$  and  $C$  are double precision arrays.

The arguments  $\ell, m, n, ka, kb, kc$  are positive integers. The coefficients  $a_j$  and  $b_j$  are assumed to be stored in  $A$  and  $B$  where

$$\begin{aligned} A(1+j \cdot ka) &= a_j \quad (j = 0, 1, \dots, \ell-1) \\ B(1+j \cdot kb) &= b_j \quad (j = 0, 1, \dots, m-1). \end{aligned}$$

IERR is a variable. When the routine is called, if  $b_0 \neq 0$  then IERR is assigned the value 0 and the first  $n$  coefficients  $c_j$  of  $p(x)/q(x)$  are stored in  $C$  where  $C(1+j \cdot kc) = c_j$  for  $j = 0, 1, \dots, n-1$ .

**Error Return.** IERR = 1 if  $b_0 = 0$ . In this case, no computation is performed.

**Remarks.** It is assumed that the array  $C$  does not overlap with the arrays  $A$  and  $B$ .

**Programmer.** A. H. Morris

## REAL POWERS OF POLYNOMIALS

If  $r$  is real and  $p(x) = \sum_{j=0}^{m-1} a_j x^j$  where  $a_0 > 0$ , then the following subroutines are available for computing the first  $n$  coefficients of the series  $p(x)^r = \sum_j b_j x^j$ .

CALL PLPWR( $r, A, ka, m, B, kb, n, IERR$ )  
CALL DPLPWR( $r, A, ka, m, B, kb, n, IERR$ )

$A$  and  $B$  are arrays. PLPWR is used if  $A$  and  $B$  are real arrays and  $r$  a real number, and DPLPWR is used if  $A$  and  $B$  are double precision arrays and  $r$  a double precision number.

The arguments  $m, n, ka, kb$  are positive integers. The coefficients  $a_j$  are assumed to be stored in  $A$  where  $A(1 + j \cdot ka) = a_j$  for  $j = 0, 1, \dots, m-1$ .  $IERR$  is a variable. When the routine is called, if  $a_0 > 0$  then  $IERR$  is assigned the value 0 and the first  $n$  coefficients  $b_j$  of  $p(x)^r$  are stored in  $B$  where  $B(1 + j \cdot kb) = b_j$  for  $j = 0, 1, \dots, n-1$ .

**Error Return.**  $IERR = 1$  if  $a_0 \leq 0$ . In this case, no computation is performed.

**Remark.** It is assumed that the arrays  $A$  and  $B$  do not overlap.

**Algorithm.** If  $q = p^r$  then  $pq' = rqp'$  where  $p'$  and  $q'$  are the derivatives of  $p$  and  $q$ . Consequently,  $b_j = \frac{1}{ja_0} \sum_{i=1}^j (ri + i - j)a_i b_{j-i}$  is used for  $j \geq 1$ . Also  $b_0 = a_0^r$ .

**Programmer.** A. H. Morris

## INVERSES OF POWER SERIES

Given an analytic function  $w = f(z) = \sum_{i \geq 1} a_i z^i$  where  $f(0) = 0$ . Then the inverse function  $z = f^{-1}(w)$  exists when  $a_1 \neq 0$ , and  $f^{-1}(w) = \sum_{i \geq 1} d_i w^i$ . The subroutines PINV and DPINV are available for obtaining the coefficients  $d_i$  when the coefficients  $a_i$  are real. DPINV is a double precision routine.

```
CALL PINV(A, D, n, WK)
CALL DPINV(A, D, n, WK)
```

If PINV is called then  $A$ ,  $D$ , and  $WK$  are real arrays. Otherwise, if DPINV is called then  $A$ ,  $D$ , and  $WK$  are double precision arrays.

It is assumed that  $n \geq 2$ .  $A$  is an array containing the coefficients  $a_1, \dots, a_n$  and  $D$  is an array of dimension  $n$ . When PINV or DPINV is called, the coefficients  $d_1, \dots, d_n$  are computed and stored in  $D$ .

$WK$  is an array of dimension  $n(n+1)/2$  or larger.  $WK$  is a work space for the routine.

**Programmer.** A. H. Morris

**Reference.** Chang, Feng-cheng, "Power Series Unification and Reversion," *Applied Math and Computation* 23 (1987), pp. 7-23.

## DERIVATIVES AND INTEGRALS OF POLYNOMIALS

Let  $f(x) = \sum_{i=0}^{n-1} a_i x^i$  be a polynomial with real coefficients  $a_i$ . The polynomial can be differentiated and integrated by the following subroutine:

**CALL MPLNMV(MO,  $x_0$ ,  $n$ ,  $A$ ,  $w$ )**

$A$  is an array containing the coefficients  $a_i$  where  $A(i) = a_{i-1}$  for  $i = 1, \dots, n$ . The argument  $x_0$  is an arbitrary real number and  $w$  is a variable. MO may have the values -1, 0, 1, 2. When MPLNMV is called  $w$  is assigned the value:

$$w = \begin{cases} \int_0^{x_0} f(x) dx & \text{if MO} = -1 \\ f(x_0) & \text{if MO} = 0 \\ f'(x_0) & \text{if MO} = 1 \\ f''(x_0) & \text{if MO} = 2 \end{cases}$$

**Programmer.** Allen V. Hershey



## EVALUATION OF CHEBYSHEV EXPANSIONS

For any complex number  $z$  and integer  $n = 0, 1, \dots$  let

$$T_0(z) = 1, \quad T_1(z) = z$$

$$T_{n+2}(z) = 2zT_{n+1}(z) - T_n(z).$$

Then  $T_n(z)$  is a polynomial of degree  $n$  having the leading coefficient  $2^{n-1}$  when  $n \geq 1$ . Also  $T_n(t) = \cos(n\theta)$  when  $t = \cos \theta (0 \leq \theta \leq \pi)$ , so that  $|T_n(t)| \leq 1$  for real  $t$  where  $|t| \leq 1$ . The polynomials  $T_n(z)$  are called the *Chebyshev polynomials (of the first kind)*. If  $f(x) = a_0/2 + \sum_{i=1}^{n-1} a_i T_i(x)$  where  $a_i$  is real, then the following functions are available for computing  $f(x)$  when  $x$  is real.

### CSEVL( $x, A, n$ )

It is assumed that  $n \geq 1$  and that  $A$  is an array containing  $a_0, a_1, \dots, a_{n-1}$  where  $A(i) = a_{i-1} (i = 1, \dots, n)$ . Then  $\text{CSEVL}(x, A, n) = f(x)$  for any real  $x$ .

**Programmer.** A. H. Morris

### DCSEVL( $x, A, n$ )

It is assumed  $n \geq 1$  and that  $A$  is a double precision array containing  $a_0, a_1, \dots, a_{n-1}$  where  $A(i) = a_{i-1} (i = 1, \dots, n)$ . Then for any double precision value  $x$ ,  $\text{DCSEVL}(x, A, n)$  is the double precision value of  $f(x)$ .

**Remark.** DCSEVL must be declared in the calling program to be of type DOUBLE PRECISION.

**Programmer.** A. H. Morris

## LAGRANGE POLYNOMIALS

Let  $a_1, \dots, a_n$  be  $n$  distinct real numbers. Then the  $i^{\text{th}}$  *Lagrange polynomial* is defined by

$$\phi_i(x) = \frac{(x - a_1)(x - a_2) \cdots (x - a_{i-1})(x - a_{i+1}) \cdots (x - a_n)}{(a_i - a_1)(a_i - a_2) \cdots (a_i - a_{i-1})(a_i - a_{i+1}) \cdots (a_i - a_n)}$$

for  $i = 1, 2, \dots, n$ . The Lagrange polynomials have the property that  $\phi_i(a_i) = 1$ ,  $\phi_i(a_j) = 0$  for  $j \neq i$ , and  $p(x) = \sum_{i=1}^n p(a_i)\phi_i(x)$  for any polynomial of degree  $n - 1$ . For convenience,

$$d_i = (a_i - a_1)(a_i - a_2) \cdots (a_i - a_{i-1})(a_i - a_{i+1}) \cdots (a_i - a_n)$$

is called the *normalization divisor* of  $\phi_i(x)$ . The following subroutines are available for computing the Lagrange polynomials and their normalization divisors.

### CALL LGRNGN(A, n, D)

$A$  and  $D$  are arrays of dimension  $n$ . The arguments  $a_1, \dots, a_n$  are given in the array  $A$ . The normalization divisors  $d_1, \dots, d_n$  are computed by the routine and stored in  $D$ .

**Programmer.** Allen V. Hershey

### CALL LGRNGV(MO, n, x<sub>0</sub>, A, D, F, DF, DDF)

$A$  and  $D$  are arrays of dimension  $n$ . The arguments  $a_1, \dots, a_n$  are given in  $A$  and the normalization divisors  $d_1, \dots, d_n$  are given in  $D$ . The argument  $x_0$  is an arbitrary real number and  $F, DF, DDF$  are arrays of dimension  $n$ .

The argument  $MO$  may take the values 0, 1, 2. If  $MO = 0$  then the polynomials  $\phi_i(x)$  are evaluated at  $x_0$  and the values  $\phi_i(x_0)$  stored in  $F$ . If  $MO = 1$  then the function  $\phi_i(x)$  and its derivative  $\phi'_i(x)$  are computed at  $x_0$ . In this case,  $\phi_i(x_0)$  is stored in  $F(i)$  and  $\phi'_i(x_0)$  is stored in  $DF(i)$  for  $i = 1, \dots, n$ . Similarly, if  $MO = 2$  then the function  $\phi_i(x)$  and its first and second derivatives are computed at  $x_0$ . The values  $\phi_i(x_0)$  are stored in  $F$ , the first derivatives are stored in  $DF$ , and the second derivatives are stored in  $DDF$ .

**Note.** If  $MO = 0$  then  $DF$  and  $DDF$  are ignored by the routine. Similarly, if  $MO = 1$  then  $DDF$  is ignored.

**Programmer.** Allen V. Hershey

### CALL LGRNGX(A, n, C)

$A$  is an array of dimension  $n$  and  $C$  an array of dimension  $n \times (n + 1)$ . The arguments  $a_1, \dots, a_n$  are given in  $A$ . The purpose of the routine is to compute the coefficients  $c_{i,j}$  of the Lagrange polynomials

$$\phi_j(x) = \sum_{k=0}^{n-1} c_{k+1,j} x^k.$$

When LGRNGX is called, the coefficients of  $\phi_j(x)$  are stored in the  $j^{\text{th}}$  column of  $C$  for  $j \leq n$ . Also, the first  $n$  coefficients of the polynomial

$$g(x) = (x - a_1) \cdots (x - a_n)$$

are stored in the  $(n + 1)^{\text{st}}$  column of  $C$ .

**Programmer.** Allen V. Hershey

## ORTHOGONAL POLYNOMIALS ON FINITE SETS

Let  $u_1, \dots, u_n$  be  $n$  distinct real numbers. For any real-valued functions  $f, g$  defined on the points  $u_i$  let  $(f, g) = \sum_{i=1}^n f(u_i)g(u_i)$ . Then  $(f, g)$  is an inner product when  $f$  and  $g$  are regarded as functions defined only on  $u_i$ . Thus, an orthonormal set of polynomials  $\{\phi_0, \phi_1, \dots, \phi_{n-1}\}$  exists where the degree of  $\phi_j$  is  $j$  for  $j < n$ . The polynomials  $\phi_j$  are defined recursively by

$$\phi_{j+1}(u) = \frac{1}{a_j} [(u - b_j)\phi_j(u) - a_{j-1}\phi_{j-1}(u)]$$

where  $a_j = (\phi_{j+1}, u\phi_j)$  and  $b_j = (\phi_j, u\phi_j)$ . Here it is assumed that  $\phi_{-1} = a_{-1} = 0$  and  $\phi_0(u) = 1/\sqrt{n}$ . The following subroutines are available for computing these polynomials.

### CALL ORTHOS( $U, m, P, n, R$ )

$U$  is an array containing the values  $u_1, \dots, u_n$  and  $m$  is an integer such that  $1 \leq m \leq n$ .  $P$  is an array of dimension  $n \times m$  and  $R$  an array of dimension  $2m - 2$ . When ORTHOS is called,  $\phi_{j-1}(u_i)$  is computed and stored in  $P(i, j)$  for  $i \leq n$  and  $j \leq m$ . Also the coefficients  $a_0, b_0, a_1, b_1, \dots, a_{m-2}, b_{m-2}$  are stored in  $R$ .

**Programmer.** Allen V. Hershey.

### CALL ORTHOV( $MO, n, u, R, m, F, DF, DDF$ )

The argument  $u$  is a real number and  $m$  an integer such that  $1 \leq m \leq n$ .  $R$  is an array containing the coefficients  $a_0, b_0, a_1, b_1, \dots, a_{m-2}, b_{m-2}$  and  $F, DF, DDF$  are arrays of dimension  $m$ .

$MO$  may take the values 0, 1, 2. If  $MO = 0$  then  $\phi_0, \phi_1, \dots, \phi_{m-1}$  are evaluated at  $u$  and the values  $\phi_{j-1}(u)$  stored in  $F$ . If  $MO = 1$  then  $\phi_{j-1}$  and its derivative  $\phi'_{j-1}$  are computed at  $u$ . In this case,  $\phi_{j-1}(u)$  is stored in  $F(j)$  and  $\phi'_{j-1}(u)$  is stored in  $DF(j)$  for  $j = 1, \dots, m$ . Similarly, if  $MO = 2$  then  $\phi_{j-1}$  and its first and second derivatives are evaluated at  $u$ . The values  $\phi_{j-1}(u)$  are stored in  $F$ , the first derivatives are stored in  $DF$ , and the second derivatives are stored in  $DDF$ .

**Note.** If  $MO = 0$  then  $DF$  and  $DDF$  are ignored by the routine. Similarly, if  $MO = 1$  then  $DDF$  is ignored.

**Programmer.** Allen V. Hershey.

### CALL ORTHOX( $n, R, m, C$ )

The argument  $m$  is an integer such that  $1 \leq m \leq n$ .  $R$  is an array containing the coefficients  $a_0, b_0, a_1, b_1, \dots, a_{m-2}, b_{m-2}$  and  $C$  an array of dimension  $m \times m$ . The purpose of the routine is to compute the coefficients  $c_{ij}$  of each polynomial

$$\phi_{j-1}(u) = \sum_{k=0}^{m-1} c_{k+1,j} u^k \quad (j = 1, \dots, m).$$

When ORTHOX is called, the coefficients of  $\phi_{j-1}$  are stored in the  $j^{\text{th}}$  column of  $C$ .

**Programmer.** Allen V. Hershey

## ZEROS OF CONTINUOUS FUNCTIONS

Let  $F(x)$  be a continuous real-valued function defined for  $a \leq x \leq b$ , and assume that  $F(a)$  and  $F(b)$  have opposite signs. Then the following function is available for finding a point  $x$  in the interval  $[a, b]$  for which  $F(x) = 0$ .

**ZEROIN**( $F, a, b, \text{AERR}, \text{RERR}$ )

ZEROIN returns a value  $x$  in the interval  $[a, b]$  for which  $F(x) = 0$ . AERR and RERR are absolute and relative tolerances that specify the desired accuracy of  $x$ . It is assumed that  $\text{AERR} \geq 0$  and  $\text{RERR} \geq 0$ . If  $x_0$  is the zero of  $F$  being approximated by  $x$ , then ZEROIN terminates when a value  $x$  is obtained for which it is estimated that  $|x - x_0| \leq \text{RERR} |x_0| + \text{AERR}$  is satisfied.

**Note.** The function  $F$  must be declared in the calling program to be of type EXTERNAL.

**Programming.** ZEROIN is a slightly modified translation of the ALGOL 60 procedure ZERO given in reference (1). The code was distributed by G. E. Forsythe, M. A. Malcolm, and C. B. Moler (University of New Mexico), and modified by A. H. Morris. The function SPMPAR is called.

### References.

- (1) Brent, Richard, *Algorithms for Minimization without Derivatives*, Prentice-Hall, 1973.
- (2) Forsythe, G.E., Malcolm, M.A., and Moler, C.B., *Computer Methods for Mathematical Computations*, Prentice-Hall, 1977.

## SOLUTION OF SYSTEMS OF NONLINEAR EQUATIONS

Let  $f_i(x) = 0$  ( $i = 1, \dots, n$ ) denote a system of  $n$  equations in  $n$  unknowns where  $x = (x_1, \dots, x_n)$ . Assume that each  $f_i(x)$  is differentiable and that an initial guess  $a = (a_1, \dots, a_n)$  to a solution of the equations is given. Then the following subroutine is available for solving the equations to within a specified tolerance.

**CALL HBRD( $F, n, X, FVEC, EPS, TOL, INFO, WK, \ell$ )**

$X$  and  $FVEC$  are arrays of dimension  $n$  or larger. On input  $X$  contains the starting point  $a = (a_1, \dots, a_n)$ . When HBRD terminates,  $X$  contains the final estimate  $x = (x_1, \dots, x_n)$  of the solution vector and  $FVEC$  contains the values of the functions  $f_1, \dots, f_n$  at the output point in  $X$ .

The argument  $F$  is the name of a user defined subroutine that has the format:

**CALL  $F(n, X, FVEC, IFLAG)$**

Here  $X$  and  $FVEC$  are arrays of dimension  $n$  and  $IFLAG$  is an integer variable. The array  $X$  contains a point  $x = (x_1, \dots, x_n)$ . Normally  $F$  will evaluate the functions  $f_1, \dots, f_n$  at this point and store the results in  $FVEC$ . However, if  $x$  does not lie in the domain of  $f_1, \dots, f_n$  then this cannot be done. In this case, the argument  $IFLAG$  (which will have been assigned a nonnegative value by HBRD) should be reset by  $F$  to a negative value. This will signal HBRD to terminate.  $F$  must be declared in the calling program to be of type EXTERNAL.

$EPS$  is an input argument which specifies the relative accuracy of  $F$ . If it is estimated that the subroutine  $F$  produces results accurate to  $k$  significant decimal digits then one may set  $EPS = 10^{-k}$ . It is required that  $EPS \geq 0$ . If  $EPS = 0$  then it is assumed that  $F$  produces results accurate to machine precision.

$TOL$  is an input argument which specifies the desired accuracy of the solution. The Euclidean norm  $\|x\| = \sqrt{\sum_i x_i^2}$  is employed. If  $\bar{x}$  denotes an actual solution of the equations, then HBRD terminates when an iterate  $x$  is generated for which it is estimated that  $\|x - \bar{x}\| \leq TOL \cdot \|x\|$  is satisfied. It is required that  $TOL \geq 0$ . In order for the convergence test to work properly, it is recommended that  $TOL$  always be smaller than  $10^{-5}$ .

$WK$  is an array of dimension  $\ell$  that is used for a work space. It is assumed that the argument  $\ell$  is greater than or equal to  $n(3n + 13)/2$ .

$INFO$  is an integer variable that reports the status of the results. When HBRD terminates,  $INFO$  has one of the following values:

$INFO < 0$  This occurs when the user terminates the execution of HBRD by resetting the argument  $IFLAG$  in the subroutine  $F$  to a negative value. Then  $INFO =$  the negative value of  $IFLAG$ .

$INFO = 0$  (Input Error)  $n < 1$ ,  $EPS < 0$ ,  $TOL < 0$ , or  $\ell < n(3n + 13)/2$ .

$INFO = 1$  A solution having the desired accuracy was obtained.

INFO = 2 The number of calls to the subroutine  $F$  has reached or exceeded  $200(n+1)$ .

INFO = 3 TOL is too small. No further improvement in the accuracy of  $x$  is possible.

INFO = 4 The routine is making very poor progress.

When HBRD terminates, if  $\text{INFO} \neq 0$  then  $X$  contains the final iterate that was generated. Also, if  $\text{INFO} \geq 1$  then FVEC contains the values of the functions  $f_1, \dots, f_n$  at this iterate. If  $\text{INFO} = 2$  then it may be helpful to continue the procedure by recalling HBRD with the current point in  $X$  as the new starting point. However, this is not advisable when  $\text{INFO} = 4$ . This setting can arise when  $x = 0$  is a solution or when entrapment occurs. HBRD searches for a solution to the equations by minimizing  $\sum_i f_i(x)^2$ . In doing this, it can become trapped in a region where the minimum does not correspond to a solution of the equations. This is what occurs when the equations have no solution. When entrapment occurs and the equations are known to have a solution, then it is recommended that the user try a different starting point.

**Scaling.** If the convergence criterion  $\|x - \bar{x}\| \leq \text{TOL} \cdot \|\bar{x}\|$  is satisfied and  $\text{TOL} = 10^{-k}$ , then the larger components of the final iterate  $x$  may be accurate to  $k$  significant digits but not the smaller components. For example, if  $\text{TOL} = 10^{-5}$  and  $x = (1.2, .34\text{E-}4)$ , then 1.2 may be accurate to 5 significant digits while .34E-4 is accurate to only 1 significant digit. If it is suspected that the smaller components do not have acceptable accuracy, then it is recommended that the variables in the original problem be rescaled and the problem rerun.

**Algorithm.** A modified form of the hybrid Powell procedure is employed.

**Programming.** HBRD is a slightly modified version of the MINPACK-1 subroutine HYBRD1. The MINPACK-1 subroutines HYBRD, ENORM, DOLGEG, FDJAC1, QFORM, QRFAC, RIMPYQ, and R1UPDT are employed. The subroutines were written by Jorge J. More, Burton S. Garbow, and Kenneth E. Hillstom (Argonne National Laboratory). The function SPMPAR is also used.

#### References.

- (1) More, J. J., Garbow, B.S., and Hillstom, K. E., *User Guide for MINPACK-1*, Argonne National Laboratory Report ANL-80-74, Argonne, Illinois, 1980.
- (2) Powell, M. J. D., "A Hybrid Method for Nonlinear Equations," *Numerical Methods for Nonlinear Algebraic Equations*, P. Rabinowitz (ed.), Gordon and Breach, London, 1970.



## SOLUTIONS OF QUADRATIC, CUBIC, AND QUARTIC EQUATIONS

Given a polynomial  $a_0 + a_1z + \cdots + a_nz^n$  with real coefficients where  $a_n \neq 0$  and  $n = 2, 3$  or  $4$ . The following subroutines are available for computing the roots  $z_1, \dots, z_n$  of the polynomial.

```
CALL QDCRT(A, Z)
CALL CBCRT(A, Z)
CALL QTCRT(A, Z)
```

It is assumed that  $A$  is a real array and  $Z$  a complex array. QDCRT is used if  $n = 2$ , CBCRT is used if  $n = 3$ , and QTCRT is used if  $n = 4$ .  $A$  is the array of coefficients where  $A(k) = a_{k-1}$  for  $k = 1, 2, \dots, n+1$ , and  $Z$  is an array of dimension  $n$ . When the appropriate subroutine is called, the roots  $z_1, \dots, z_n$  are stored in  $Z$ . The real roots precede the complex roots. The real roots are ordered so that  $|z_j| \leq |z_{j+1}|$ . The complex roots are unordered except that complex conjugate pairs of roots appear consecutively with the root having the positive imaginary part being first.

**Programming.** QTCRT calls the subroutines CBCRT and AORD, and CBCRT calls the subroutine QDCRT and function CBRT. The routines were written by A. H. Morris and CBCRT was modified by Wm. Davis (NSWC). The function SPMPAR is also used.

```
CALL DQDCRT(A, ZR, ZI)
CALL DCBCRT(A, ZR, ZI)
CALL DQTCRT(A, ZR, ZI)
```

It is assumed that  $A, ZR$ , and  $ZI$  are double precision arrays. DQDCRT is used if  $n = 2$ , DCBCRT is used if  $n = 3$ , and DQTCRT is used if  $n = 4$ .  $A$  is the array of coefficients where  $A(k) = a_{k-1}$  for  $k = 1, 2, \dots, n+1$ , and  $ZR$  and  $ZI$  are arrays of dimension  $n$ . When the appropriate subroutine is called, the real parts of the roots  $z_1, \dots, z_n$  are stored in  $ZR$  and the imaginary parts are stored in  $ZI$ . The real roots precede the complex roots. The real roots are ordered so that  $|z_j| \leq |z_{j+1}|$ . The complex roots are unordered except that complex conjugate pairs of roots appear consecutively with the root having the positive imaginary part being first.

**Programming.** DQTCRT calls the routines DAORD, DCSQRT, and DCBCRT. DCBCRT calls the subroutine DQDCRT and function DCBRT. The routines were written by A. H. Morris. The function DPMPAR is also used.

## DOUBLE PRECISION ROOTS OF POLYNOMIALS

Given a polynomial  $a_0 + a_1z + \dots + a_nz^n$  of degree  $n \geq 1$ . The subroutines DRPOLY and DCPOLY are available for obtaining the roots  $z_1, \dots, z_n$  of the polynomial. DRPOLY may be used if the coefficients  $a_j$  are real, and DCPOLY is applicable if the coefficients are complex. These subroutines perform the calculations in double precision.

### CALL DRPOLY(A,n,ZR,ZI,NUM,WK,DWK)

A is a double precision array containing the coefficients where  $A(j) = a_{n-j+1}$  for  $j = 1, \dots, n+1$ . ZR and ZI are double precision arrays of dimension  $n$  or larger, and NUM is an integer variable. When DRPOLY is called, if no errors are detected then NUM = the number of roots that are obtained. If  $NUM \geq 1$  then the real parts of the roots are stored in ZR(j) and the imaginary parts in ZI(j) for  $j = 1, \dots, NUM$ . The roots are unordered except that complex conjugate pairs of roots appear consecutively with the positive imaginary part being first.

WK is a real array of dimension  $n+1$  or larger, and DWK is a double precision array of dimension  $6(n+1)$  or larger. WK and DWK are work spaces for the routine.

**Error Return.** NUM = -1 if  $n < 1$  or  $a_n = 0$ .

**Programming.** DRPOLY employs the subroutines DRPLY1,FXSHFR,QUADIT,REALIT,CALCSC,NEXTK,NEWEST,QUADSD, and QUADPL. These routines exchange information in a labeled common block named GLOBAL. The routines were written by M. A. Jenkins (Queen's University, Kingston, Ontario), and modified by A. H. Morris. The functions SPMPAR, DPMPAR, and IPMPAR are also used.

### References.

- (1) Jenkins, M. A., "Zeros of a Real Polynomial," *ACM Trans. Math Software* 1 (1975), pp. 178-189.
- (2) Jenkins, M. A. and Traub, J. F., "A Three-Stage Algorithm for Real Polynomials using Quadratic Iterations," *SIAM J. Numerical Analysis* 7 (1970), pp. 545-566.

### CALL DCPOLY(AR,AI,n,ZR,ZI,NUM,DWK)

AR and AI are double precision arrays containing the real and imaginary parts of the coefficients where  $AR(j) = \text{Re}(a_{n-j+1})$  and  $AI(j) = \text{Im}(a_{n-j+1})$  for  $j = 1, \dots, n+1$ . ZR and ZI are double precision arrays of dimension  $n$  or larger, and NUM is an integer variable. When DCPOLY is called, if no errors are detected then NUM = the number of roots that are obtained. If  $NUM \geq 1$  then the real parts of the roots are stored in ZR(j) and the imaginary parts in ZI(j) for  $j = 1, \dots, NUM$ .

DWK is a double precision array of dimension  $10(n+1)$  or larger that is a work space for the routine.

**Error Return.**  $\text{NUM} = -1$  if  $n < 1$  or  $a_n = 0$ .

**Programming.** DCPOLY employs the routines DCPLY1, CAUCHY, NOSHFT, FXSHFT, VRSHFT, CALCT, NEXTH, POLYEV, CDIVID, and the functions SCALCP, ERREV. These routines and functions were written by M.A. Jenkins (Queen's University, Kingston, Ontario) and J. F. Traub (Bell Laboratories), and modified by A. H. Morris. The functions DPMPAR, IPMPAR, and DCPABS are also used.

#### **References.**

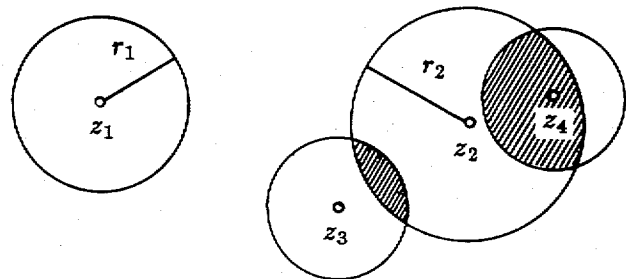
- (1) Jenkins, M. A. and Traub, J. F., "Algorithm 419, Zeros of a Complex Polynomial," *Comm. ACM* **15** (1972), pp. 97-99.
- (2) \_\_\_\_\_, "A Three-Stage Variable-Shift Iteration for Polynomial Zeros and its Relation to Generalized Rayleigh Iteration," *Numer. Math* **14** (1970), pp. 252-263.

## ACCURACY OF THE ROOTS OF A REAL POLYNOMIAL

Given a polynomial  $p(z) = a_0 + a_1z + \dots + a_nz^n$  of degree  $n \geq 1$  with real coefficients. Let  $z_1, \dots, z_n$  be approximations for the roots of  $p(z)$ . Then for each  $z_i$ , the subroutine RBND obtains the radius  $r_i$  of a disk  $D_i = \{z : |z - z_i| \leq r_i\}$  centered at  $z_i$  which contains a true zero of the polynomial. The radius  $r_i$  is an upper bound on the absolute error of the approximation  $z_i$ .

For each  $z_i$ , RBND also computes the number  $k_i$  of disks  $D_j$  (including the disk  $D_i$ ) which overlap with  $D_i$ . The value  $k_i$  is the number of roots of  $p(z)$  that are clustered near  $z_i$ . If  $k_i = 1$  then  $z_i$  approximates a simple root.

**Example.** In the figure  $k_1 = 1$ ,  $k_2 = 3$ ,  $k_3 = 2$ , and  $k_4 = 2$ .



CALL RBND( $n, A, Z, R, RERR, K, IERR$ )

$A$  is a real array containing the coefficients  $a_0, a_1, \dots, a_n$  and  $Z$  a complex array containing the approximate roots  $z_1, \dots, z_n$ . It is assumed that  $n \geq 1$  and  $A(i) = a_{i-1}$  for  $i = 1, \dots, n+1$ .

$IERR$  is an integer variable,  $R$  a real array of dimension  $n$  or larger, and  $K$  an integer array of dimension  $n$  or larger. When RBND is called, if no input errors are detected then  $IERR$  is assigned the value 0, the radii  $r_1, \dots, r_n$  are computed and stored in  $R$ , and the values  $k_1, \dots, k_n$  are computed and stored in  $K$ .

$RERR$  is a real array of dimension  $n$  or larger. If  $z_i = 0$  then  $RERR(i)$  is set to  $-1$  by the routine. Otherwise, if  $z_i \neq 0$  then  $RERR(i)$  = the estimated maximum relative error of  $z_i$ .

**Error Return.**  $IERR = 1$  when  $n < 1$  and  $IERR = 2$  when  $a_n = 0$ . In these cases no computation is performed.

**Programming.** RBND employs the functions CPABS and SPMPAR. RBND was written by Carl B. Bailey and modified by William R. Gavin (Sandia Laboratories). The format of the routine was modified by A.H. Morris.

## COPYING VECTORS

A copy of a vector  $X = (x_1, \dots, x_n)$  can be made and stored in the array  $Y$  by the following subroutines:

```
CALL SCOPY( $n, X, kx, Y, ky$ )  
CALL DCOPY( $n, X, kx, Y, ky$ )  
CALL CCOPY( $n, X, kx, Y, ky$ )
```

The argument  $kx$  is an integer which specifies the interval between successive components  $x_i$  of the vector  $X$ . If  $kx \geq 0$  then it is assumed that  $x_i$  is stored in  $X(1+(i-1)kx)$  for  $i = 1, \dots, n$ . Otherwise, if  $kx < 0$  then it is assumed that  $x_i$  is stored in  $X(1+(n-i)|kx|)$ . Similarly, the argument  $ky$  specifies the spacing of the components of  $Y$ .

SCOPY is used if  $X$  and  $Y$  are real arrays, DCOPY is used if  $X$  and  $Y$  are double precision arrays, and CCOPY is used if  $X$  and  $Y$  are complex arrays. When any of these routines is called, if  $n \leq 0$  then the routine immediately terminates. Otherwise, if  $n \geq 1$  then the components  $x_i$  of  $X$  are stored in  $Y$  according to the spacing specified by the  $ky$  parameter.

**Programming.** These routines are part of the BLAS package of basic linear algebra subroutines designed by C. L. Lawson, R. J. Hanson, D. R. Kincaid, and F. T. Krogh. The routines were written by Jack Dongarra (Argonne National Laboratory).

**Reference.** Lawson, C. L., Hanson, R. J., Kincaid, D. R., and Krogh, F. T., *Basic Linear Algebra Subprograms for FORTRAN Usage*. Report SAND 77-0898, Sandia Laboratories, Albuquerque, New Mexico, 1977.

## INTERCHANGING VECTORS

The components of two vectors  $X = (x_1, \dots, x_n)$  and  $Y = (y_1, \dots, y_n)$  can be interchanged by the following subroutines:

```
CALL SSWAP( $n, X, kx, Y, ky$ )  
CALL DSWAP( $n, X, kx, Y, ky$ )  
CALL CSWAP( $n, X, kx, Y, ky$ )
```

The argument  $kx$  is an integer which specifies the interval between successive components  $x_i$  of the vector  $X$ . If  $kx \geq 0$  then it is assumed that  $x_i$  is stored in  $X(1 + (i-1)kx)$  for  $i = 1, \dots, n$ . Otherwise, if  $kx < 0$  then it is assumed that  $x_i$  is stored in  $X(1 + (n-i)|kx|)$ . Similarly, the argument  $ky$  specifies the spacing of the components of  $Y$ .

SSWAP is used if  $X$  and  $Y$  are real arrays, DSWAP is used if  $X$  and  $Y$  are double precision arrays, and CSWAP is used if  $X$  and  $Y$  are complex arrays. When any of these routines is called, if  $n \leq 0$  then the routine immediately terminates. Otherwise, if  $n \geq 1$  then the components  $x_i$  and  $y_i$  are interchanged for  $i = 1, \dots, n$ . Thus, when the routine terminates  $X = (y_1, \dots, y_n)$  and  $Y = (x_1, \dots, x_n)$ .

**Programming.** These routines are part of the BLAS package of basic linear algebra subroutines designed by C. L. Lawson, R. J. Hanson, D. R. Kincaid, and F. T. Krogh. The routines were written by Jack Dongarra (Argonne National Laboratory).

**Reference.** Lawson, C. L., Hanson, R. J., Kincaid, D. R., and Krogh, F. T., *Basic Linear Algebra Subprograms for FORTRAN Usage*. Report SAND 77-0898, Sandia Laboratories, Albuquerque, New Mexico, 1977.

## PLANAR ROTATION OF VECTORS

Let  $X = (x_1, \dots, x_n)$  and  $Y = (y_1, \dots, y_n)$  be vectors and  $c$  and  $s$  be real numbers such that  $c^2 + s^2 = 1$ .  $X$  and  $Y$  can be replaced with  $cX + sY$  and  $-sX + cY$  by the following subroutines:

```
CALL SROT( $n, X, kx, Y, ky, c, s$ )  
CALL DROT( $n, X, kx, Y, ky, c, s$ )  
CALL CSROT( $n, X, kx, Y, ky, c, s$ )
```

The argument  $kx$  is an integer which specifies the interval between successive components  $x_i$  of the vector  $X$ . If  $kx \geq 0$  then it is assumed that  $x_i$  is stored in  $X(1 + (i-1)kx)$  for  $i = 1, \dots, n$ . Otherwise, if  $kx < 0$  then it is assumed that  $x_i$  is stored in  $X(1 + (n-i)|kx|)$ . Similarly, the argument  $ky$  specifies the spacing of the components of  $Y$ .

SROT is used if  $X$  and  $Y$  are real arrays, DROT is used if  $X$  and  $Y$  are double precision arrays, and CSROT is used if  $X$  and  $Y$  are complex arrays. The arguments  $c$  and  $s$  are real numbers when SROT and CSROT are used, and are double precision numbers when DROT is used. When any of these routines is called, if  $n \leq 0$  then the routine immediately terminates. Otherwise, if  $n \geq 1$  then the components  $x_i$  and  $y_i$  are replaced with  $cx_i + sy_i$  and  $-sx_i + cy_i$  for  $i = 1, \dots, n$ .

**Programming.** These routines are part of the BLAS package of basic linear algebra subroutines designed by C. L. Lawson, R. J. Hanson, D. R. Kincaid, and F. T. Krogh. The routines were written by Jack Dongarra (Argonne National Laboratory).

**Reference.** Lawson, C. L., Hanson, R. J., Kincaid, D. R., and Krogh, F. T., *Basic Linear Algebra Subprograms for FORTRAN Usage*. Report SAND 77-0898, Sandia Laboratories, Albuquerque, New Mexico, 1977.

## DOT PRODUCTS OF VECTORS

The following functions are available for computing the sums  $\sum_i x_i y_i$  and  $\sum_i \bar{x}_i y_i$  where  $X = (x_1, \dots, x_n)$  and  $Y = (y_1, \dots, y_n)$  are real or complex vectors.

**SDOT**( $n, X, kx, Y, ky$ )  
**DDOT**( $n, X, kx, Y, ky$ )  
**CDOTC**( $n, X, kx, Y, ky$ )  
**CDOTU**( $n, X, kx, Y, ky$ )

The argument  $kx$  is an integer which specifies the interval between successive components  $x_i$  of the vector  $X$ . If  $kx \geq 0$  then it is assumed that  $x_i$  is stored in  $X(1+(i-1)kx)$  for  $i = 1, \dots, n$ . Otherwise, if  $kx < 0$  then it is assumed that  $x_i$  is stored in  $X(1+(n-i)|kx|)$ . Similarly, the argument  $ky$  specifies the spacing of the components of  $Y$ .

SDOT is used if  $X$  and  $Y$  are real arrays, and DDOT is used if  $X$  and  $Y$  are double precision arrays. SDOT is a real function and DDOT is a double precision function. When either of these two functions is called, if  $n \leq 0$  then the function is assigned the value 0. Otherwise, if  $n \geq 1$  then the function is assigned the value  $\sum_{i=1}^n x_i y_i$ .

CDOTC and CDOTU are used if  $X$  and  $Y$  are complex arrays. CDOTC and CDOTU are complex functions. When either of these two functions is called, if  $n \leq 0$  then the function is assigned the value 0. Otherwise, if  $n \geq 1$  then CDOTC( $n, X, kx, Y, ky$ ) is assigned the value  $\sum_{i=1}^n \bar{x}_i y_i$  and CDOTU( $n, X, kx, Y, ky$ ) is assigned the value  $\sum_{i=1}^n x_i y_i$ .

**Remark.** DDOT must be declared in the calling program to be of type DOUBLE PRECISION, and CDOTC and CDOTU must be declared to be of type COMPLEX.

**Programming.** These functions are part of the BLAS package of basic linear algebra subroutines designed by C. L. Lawson, R. J. Hanson, D. R. Kincaid, and F. T. Krogh. The functions were written by Jack Dongarra (Argonne National Laboratory).

**Reference.** Lawson, C. L., Hanson, R. J., Kincaid, D. R., and Krogh, F. T., *Basic Linear Algebra Subprograms for FORTRAN Usage*. Report SAND 77-0898, Sandia Laboratories, Albuquerque, New Mexico, 1977.



## SCALING VECTORS

If  $a$  is a real or complex number and  $X = (x_1, \dots, x_n)$  a vector, then the vector  $X$  can be replaced with the vector  $aX$  by the following subroutines:

```
CALL SSCAL( $n, a, X, kx$ )  
CALL DSCAL( $n, a, X, kx$ )  
CALL CSCAL( $n, a, X, kx$ )  
CALL CSSCAL( $n, a, X, kx$ )
```

The argument  $kx$  is a positive integer. It is assumed that the component  $x_i$  is stored in  $X(1 + (i - 1)kx)$  for  $i = 1, \dots, n$ .

SSCAL is used if  $a$  is a real number and  $X$  a real array, DSCAL is used if  $a$  is a double precision number and  $X$  a double precision array, CSCAL is used if  $a$  is a complex number and  $X$  a complex array, and CSSCAL is used if  $a$  is a real number and  $X$  a complex array. When any of these routines is called, if  $n \leq 0$  then the routine immediately terminates. Otherwise, if  $n \geq 1$  then each  $x_i$  is replaced with  $ax_i$ . Thus when the routine terminates  $X = (ax_1, \dots, ax_n)$ .

**Programming.** These routines are part of the BLAS package of basic linear algebra subroutines designed by C. L. Lawson, R. J. Hanson, D. R. Kincaid, and F. T. Krogh. The routines were written by Jack Dongarra (Argonne National Laboratory).

**Reference.** Lawson, C. L., Hanson, R. J., Kincaid, D. R., and Krogh, F. T., *Basic Linear Algebra Subprograms for FORTRAN Usage*. Report SAND 77-0898, Sandia Laboratories, Albuquerque, New Mexico, 1977.

## VECTOR ADDITION

If  $a$  is a real or complex number and  $X = (x_1, \dots, x_n)$  a vector, then any vector  $Y = (y_1, \dots, y_n)$  can be replaced with the vector  $aX + Y$  by the following subroutines:

```
CALL SAXPY( $n, a, X, kx, Y, ky$ )  
CALL DAXPY( $n, a, X, kx, Y, ky$ )  
CALL CAXPY( $n, a, X, kx, Y, ky$ )
```

The argument  $kx$  is an integer which specifies the interval between successive components  $x_i$  of the vector  $X$ . If  $kx \geq 0$  then it is assumed that  $x_i$  is stored in  $X(1+(i-1)kx)$  for  $i = 1, \dots, n$ . Otherwise, if  $kx < 0$  then it is assumed that  $x_i$  is stored in  $X(1+(n-i)|kx|)$ . Similarly, the argument  $ky$  specifies the spacing of the components of the vector  $Y$ .

SAXPY is used if  $a$  is a real number and  $X, Y$  are real arrays, DAXPY is used if  $a$  is a double precision number and  $X, Y$  are double precision arrays, and CAXPY is used if  $a$  is a complex number and  $X, Y$  are complex arrays. When any of these routines is called, if  $n \leq 0$  or  $a = 0$  then the routine immediately terminates. Otherwise, if  $n \geq 1$  then  $y_i$  is replaced with  $ax_i + y_i$  for  $i = 1, \dots, n$ . Thus when the routine terminates  $Y = (ax_1 + y_1, \dots, ax_n + y_n)$ .

**Programming.** These routines are part of the BLAS package of basic linear algebra subroutines designed by C. L. Lawson, R. J. Hanson, D. R. Kincaid, and F. T. Krogh. The routines were written by Jack Dongarra (Argonne National Laboratory).

**Reference.** Lawson, C. L., Hanson, R. J., Kincaid, D. R., and Krogh, F. T., *Basic Linear Algebra Subprograms for FORTRAN Usage*. Report SAND 77-0898, Sandia Laboratories, Albuquerque, New Mexico, 1977.

## L<sub>1</sub> NORM OF A VECTOR

The following functions are available for computing the  $L_1$  norm of a real vector or a modified  $L_1$  norm of a complex vector.

**SASUM**( $n, X, kx$ )  
**DASUM**( $n, X, kx$ )  
**SCASUM**( $n, X, kx$ )

$X = (x_1, \dots, x_n)$  is a vector and  $kx$  a positive integer. It is assumed that  $x_i$  is stored in  $X(1 + (i - 1)kx)$  for  $i = 1, \dots, n$ .

SASUM is used if  $X$  is a real array and DASUM is used if  $X$  is a double precision array. SASUM is a real function and DASUM a double precision function. When either of these functions is called, if  $n \leq 0$  then the function is assigned the value 0. Otherwise, if  $n \geq 1$  then the function is assigned the value  $\sum_{i=1}^n |x_i|$ .

SCASUM is used if  $X$  is a complex array. SCASUM is a real function. When SCASUM is called, if  $n \leq 0$  then the function is assigned the value 0. Otherwise, if  $n \geq 1$  then SCASUM( $n, X, kx$ ) is assigned the value  $\sum_{i=1}^n [|Re(x_i)| + |Im(x_i)|]$ .

### Remarks.

- (1) DASUM must be declared in the calling program to be of type DOUBLE PRECISION.
- (2) SCASUM( $n, X, kx$ ) is the norm of the complex vector  $X = (x_1, \dots, x_n)$  when  $X$  is regarded as a real vector of dimension  $2n$ . This norm is frequently preferred over the standard  $L_1$  norm  $\sum_{i=1}^n |x_i|$  since it takes less time to compute.

**Programming.** These functions are part of the BLAS package of basic linear algebra subroutines designed by C. L. Lawson, R. J. Hanson, D. R. Kincaid, and F. T. Krogh. The functions were written by Jack Dongarra (Argonne National Laboratory).

**Reference.** Lawson, C. L., Hanson, R. J., Kincaid, D. R., and Krogh, F. T., *Basic Linear Algebra Subprograms for FORTRAN Usage*. Report SAND 77-0898, Sandia Laboratories, Albuquerque, New Mexico, 1977.

## L<sub>2</sub> NORM OF A VECTOR

The following functions are available for computing the  $L_2$  norm of a real or complex vector.

SNRM2( $n, X, kx$ )  
DNRM2( $n, X, kx$ )  
SCNRM2( $n, X, kx$ )

$X = (x_1, \dots, x_n)$  is a vector and  $kx$  a positive integer. It is assumed that  $x_i$  is stored in  $X(1 + (i - 1)kx)$  for  $i = 1, \dots, n$ .

SNRM2 is used if  $X$  is a real array, DNRM2 is used if  $X$  is a double precision array, and SCNRM2 is used if  $X$  is a complex array. SNRM2 and SCNRM2 are real functions and DNRM2 a double precision function. When any of these functions is called, if  $n \leq 0$  then the function is assigned the value 0. Otherwise, if  $n \geq 1$  then the function is assigned the value  $\left[ \sum_{i=1}^n |x_i|^2 \right]^{1/2}$ .

**Remark.** The function DNRM2 must be declared in the calling program to be of type DOUBLE PRECISION.

**Programming.** These functions are part of the BLAS package of basic linear algebra subroutines designed by C. L. Lawson, R. J. Hanson, D. R. Kincaid, and F. T. Krogh. The functions were written by Charles Lawson (Jet Propulsion Laboratory).

**Reference.** Lawson, C. L., Hanson, R. J., Kincaid, D. R., and Krogh, F. T., *Basic Linear Algebra Subprograms for FORTRAN Usage*. Report SAND 77-0898, Sandia Laboratories, Albuquerque, New Mexico, 1977.

## $L_\infty$ NORM OF A VECTOR

The following functions are available for finding the largest component  $x_i$  of a vector  $X = (x_1, \dots, x_n)$ .

**ISAMAX**( $n, X, kx$ )

**IDAMAX**( $n, X, kx$ )

**ICAMAX**( $n, X, kx$ )

The argument  $kx$  is a positive integer. It is assumed that the component  $x_i$  is stored in  $X(1 + (i - 1)kx)$  for  $i = 1, \dots, n$ .

ISAMAX is used if  $X$  is a real array and IDAMAX is used if  $X$  is a double precision array. ISAMAX and IDAMAX are integer functions. When either of these functions is called, if  $n \leq 0$  then the function is assigned the value 0. Otherwise, if  $n \geq 1$  then the function is assigned the value  $i$  where  $i$  is the smallest index such that  $|x_i| = \max\{|x_j| : j = 1, \dots, n\}$ .

ICAMAX is also an integer function. It is used when  $X$  is a complex array. If  $n \leq 0$  then ICAMAX( $n, X, kx$ ) is assigned the value 0. Otherwise, if  $n \geq 1$  then the function is assigned the value  $i$  where  $i$  is the smallest index such that  $|\operatorname{Re}(x_i)| + |\operatorname{Im}(x_i)| = \max\{|\operatorname{Re}(x_j)| + |\operatorname{Im}(x_j)| : j = 1, \dots, n\}$ .

**Note.** The mapping  $X \rightarrow \max\{|\operatorname{Re}(x_j)| + |\operatorname{Im}(x_j)| : j = 1, \dots, n\}$  is the  $L_\infty$  norm of the complex vector  $X = (x_1, \dots, x_n)$  when  $X$  is regarded as a real  $n \times 2$  matrix. This norm is frequently preferred over the standard  $L_\infty$  norm  $\max\{|x_j| : j = 1, \dots, n\}$  since it takes less time to compute.

**Programming.** These functions are part of the BLAS package of basic linear algebra subroutines designed by C. L. Lawson, R. J. Hanson, D. R. Kincaid, and F. T. Krogh. The functions were written by Jack Dongarra (Argonne National Laboratory).

**Reference.** Lawson, C. L., Hanson, R. J., Kincaid, D. R., and Krogh, F. T., *Basic Linear Algebra Subprograms for FORTRAN Usage*. Report SAND 77-0898, Sandia Laboratories, Albuquerque, New Mexico, 1977.

## PACKING AND UNPACKING SYMMETRIC MATRICES

An  $n \times n$  symmetric matrix  $A = (a_{ij})$  can be represented by either its lower triangular elements

$$\begin{pmatrix} a_{11} & a_{12} & a_{13} & \dots \\ a_{21} & a_{22} & a_{23} & \dots \\ a_{31} & a_{32} & a_{33} & \dots \\ \vdots & \vdots & \vdots & \ddots \end{pmatrix}$$

or its upper triangular elements. If the lower triangular elements are used, then the packed form for the matrix is an array of dimension  $n(n+1)/2$  containing the elements  $a_{11}a_{21}a_{22}a_{31}a_{32}a_{33} \dots$ . Similarly, if the upper triangular elements are used then the packed form for the matrix is an array containing  $a_{11}a_{12} \dots a_{1n}a_{22}a_{23} \dots a_{2n} \dots$ . *Currently the lower triangular elements are used for packing symmetric matrices.* The following subroutines are available for packing and unpacking real symmetric matrices.

**CALL MCVFS( $A, ka, n, B$ )**  
**CALL DMCVFS( $A, ka, n, B$ )**

$A$  is an  $n \times n$  symmetric matrix and  $B$  an array whose dimension is equal to or greater than  $n(n+1)/2$ . The routines store the packed form of  $A$  in  $B$ . MCVFS is used if  $A$  and  $B$  are real arrays and DMCVFS is used if  $A$  and  $B$  are double precision arrays. The input argument  $ka$  has the following value:

$ka$  = the number of rows in the dimension statement for  $A$  in the calling program  
 It is assumed that  $ka \geq n$ .

**Note.**  $A$  and  $B$  may begin in the same location.

**Programmer.** A. H. Morris

**CALL MCVSF( $A, ka, n, B$ )**  
**CALL DMCVSF( $A, ka, n, B$ )**

$B$  is an array containing the elements of a packed  $n \times n$  symmetric matrix and  $A$  is an array of dimension  $ka \times n$  where  $ka \geq n$ . The routines unpack  $B$  and store the results in  $A$ . MCVSF is used if  $A$  and  $B$  are real arrays and DMCVSF is used if  $A$  and  $B$  are double precision arrays.

**Note.**  $A$  and  $B$  may begin in the same location.

**Programmer.** A. H. Morris

## CONVERSION OF REAL MATRICES TO AND FROM DOUBLE PRECISION FORM

The subroutines MCVRD and MCVDR are available for converting real matrices to and from double precision form.

**CALL MCVRD( $m, n, A, ka, B, kb$ )**

$A$  is a real  $m \times n$  matrix and  $B$  a double precision 2-dimensional array. MCVRD stores the matrix in double precision form in the array  $B$ . The input arguments  $ka$  and  $kb$  have the following values:

$ka$  = the number of rows in the dimension statement for  $A$  in the calling program

$kb$  = the number of rows in the dimension statement for  $B$  in the calling program

It is assumed that  $ka \geq m$  and  $kb \geq m$ , and that  $B$  contains at least  $n$  columns.

**Programmer.** A. H. Morris

**CALL MCVDR( $m, n, A, ka, B, kb$ )**

$A$  is a double precision  $m \times n$  matrix and  $B$  a real 2-dimensional array. MCVDR stores the matrix in single precision form in the array  $B$ . The input arguments  $ka$  and  $kb$  have the following values:

$ka$  = the number of rows in the dimension statement for  $A$  in the calling program

$kb$  = the number of rows in the dimension statement for  $B$  in the calling program

It is assumed that  $ka \geq m$  and  $kb \geq m$ , and that  $B$  contains at least  $n$  columns.

**Programmer.** A. H. Morris

## STORAGE OF REAL MATRICES IN THE COMPLEX MATRIX FORMAT

The following subroutine is available for storing a real matrix  $A$  in the complex matrix format.

**CALL MCVRC( $m, n, A, ka, B, kb$ )**

$A$  is a real  $m \times n$  matrix and  $B$  a complex 2-dimensional array. MCVRC stores the matrix in complex form in the array  $B$ . The input arguments  $ka$  and  $kb$  have the following values:

$ka$  = the number of rows in the dimension statement for  $A$  in the calling program

$kb$  = the number of rows in the dimension statement for  $B$  in the calling program

It is assumed that  $ka \geq m$  and  $kb \geq m$ , and that  $B$  contains at least  $n$  columns.

**Programmer.** A. H. Morris



## THE REAL AND IMAGINARY PARTS OF A COMPLEX MATRIX

If  $A = (a_{ij})$  is a complex matrix then let  $\text{Re}(A) = (\text{Re}(a_{ij}))$  and  $\text{Im}(A) = (\text{Im}(a_{ij}))$  denote the real and imaginary parts of  $A$ . The following subroutines are available for obtaining  $\text{Re}(A)$  and  $\text{Im}(A)$ .

**CALL CMREAL( $m, n, A, ka, B, kb$ )**

$A$  is a complex  $m \times n$  matrix and  $B$  a real 2-dimensional array. CMREAL obtains  $\text{Re}(A)$  and stores it in  $B$ . The input arguments  $ka$  and  $kb$  have the following values:

$ka$  = the number of rows in the dimension statement for  $A$  in the calling program

$kb$  = the number of rows in the dimension statement for  $B$  in the calling program

It is assumed that  $ka \geq m$  and  $kb \geq m$ , and that  $B$  contains at least  $n$  columns.

**Programmer.** A. H. Morris

**CALL CMIMAG( $m, n, A, ka, B, kb$ )**

$A$  is a complex  $m \times n$  matrix and  $B$  a real 2-dimensional array. CMIMAG obtains  $\text{Im}(A)$  and stores it in  $B$ . The input arguments  $ka$  and  $kb$  have the following values:

$ka$  = the number of rows in the dimension statement for  $A$  in the calling program

$kb$  = the number of rows in the dimension statement for  $B$  in the calling program

It is assumed that  $ka \geq m$  and  $kb \geq m$ , and that  $B$  contains at least  $n$  columns.

**Programmer.** A. H. Morris

## COPYING MATRICES

The following subroutines are available for copying matrices.

### CALL MCOPY( $m, n, A, ka, B, kb$ )

$A$  is a real  $m \times n$  matrix and  $B$  a real 2-dimensional array. MCOPY makes a copy of the matrix  $A$  and stores it in  $B$ . The input arguments  $ka$  and  $kb$  have the following values:

$ka$  = the number of rows in the dimension statement for  $A$  in the calling program

$kb$  = the number of rows in the dimension statement for  $B$  in the calling program

It is assumed that  $ka \geq m$  and  $kb \geq m$ , and that  $B$  contains at least  $n$  columns.

Programmer. A. H. Morris

### CALL SMCOPY( $n, A, B$ )

$A$  is a real packed  $n \times n$  symmetric matrix and  $B$  a real array whose dimension is equal to or greater than  $n(n+1)/2$ . SMCOPY makes a copy of the packed symmetric matrix  $A$  and stores it in  $B$ .

Programmer. A. H. Morris

### CALL DMCOPY( $m, n, A, ka, B, kb$ )

$A$  is a double precision  $m \times n$  matrix and  $B$  a double precision 2-dimensional array. DMCOPY makes a copy of the matrix  $A$  and stores it in  $B$ . The input arguments  $ka$  and  $kb$  have the following values:

$ka$  = the number of rows in the dimension statement for  $A$  in the calling program

$kb$  = the number of rows in the dimension statement for  $B$  in the calling program

It is assumed that  $ka \geq m$  and  $kb \geq m$ , and that  $B$  contains at least  $n$  columns.

Programmer. A. H. Morris

### CALL CMCOPY( $m, n, A, ka, B, kb$ )

$A$  is a complex  $m \times n$  matrix and  $B$  a complex 2-dimensional array. CMCOPY makes a copy of the matrix  $A$  and stores it in  $B$ . The input arguments  $ka$  and  $kb$  have the following values:

$ka$  = the number of rows in the dimension statement for  $A$  in the calling program

$kb$  = the number of rows in the dimension statement for  $B$  in the calling program

It is assumed that  $ka \geq m$  and  $kb \geq m$ , and that  $B$  contains at least  $n$  columns.

Programmer. A. H. Morris

## COMPUTATION OF THE CONJUGATE OF A COMPLEX MATRIX

If  $A = (a_{ij})$  is a complex matrix then let  $\bar{A} = (\bar{a}_{ij})$  denote the conjugate of  $A$ . The following subroutine is available for computing the conjugate matrix  $\bar{A}$ .

**CALL CMCONJ( $m, n, A, ka, B, kb$ )**

$A$  is a complex  $m \times n$  matrix and  $B$  a complex 2-dimensional array. CMCONJ computes  $\bar{A}$  and stores the results in  $B$ . The input arguments  $ka$  and  $kb$  have the following values:

$ka$  = the number of rows in the dimension statement for  $A$  in the calling program

$kb$  = the number of rows in the dimension statement for  $B$  in the calling program

It is assumed that  $ka \geq m$  and  $kb \geq m$ , and that  $B$  contains at least  $n$  columns.

**Remark.**  $A$  and  $B$  may reference the same storage area when  $ka = kb$ .

**Programmer.** A. H. Morris

## TRANSPOSING MATRICES

The subroutines TPOSE, DTPOSE, CTPOSE and TIP, DTIP, CTIP are available for transposing a matrix  $A$ . TPOSE, DTPOSE, and CTPOSE are used if the results are to be stored in a separate storage area  $B$ . TIP, DTIP, and CTIP are used if the results are to be stored in  $A$  (i.e., if the transposition is to be done in place).

```
CALL TPOSE( $m, n, A, ka, B, kb$ )  
CALL DTPOSE( $m, n, A, ka, B, kb$ )  
CALL CTPOSE( $m, n, A, ka, B, kb$ )
```

TPOSE is called if  $A$  is a real matrix and  $B$  a real array, DTPOSE is called if  $A$  is a double precision matrix and  $B$  a double precision array, and CTPOSE is called if  $A$  is a complex matrix and  $B$  a complex array.

$A$  is a matrix having  $m$  rows and  $n$  columns, and  $B$  a 2-dimensional array. The routine transposes  $A$  and stores the results in  $B$ . The input arguments  $ka$  and  $kb$  have the following values:

$ka$  = the number of rows in the dimension statement for  $A$  in the calling program

$kb$  = the number of rows in the dimension statement for  $B$  in the calling program

It is assumed that  $ka \geq m$  and  $kb \geq n$ , and that  $B$  contains at least  $n$  columns.

**Programmer.** A. H. Morris

```
CALL TIP( $A, m, n, \text{MOVE}, k, \text{MDIM}$ )  
CALL DTIP( $A, m, n, \text{MOVE}, k, \text{MDIM}$ )  
CALL CTIP( $A, m, n, \text{MOVE}, k, \text{MDIM}$ )
```

TIP is called if  $A$  is a real array, DTIP is called if  $A$  is a double precision array, and CTIP is called if  $A$  is a complex array.

$A$  is an array of dimension  $mn$  which contains an  $m \times n$  matrix to be transposed. The routine transposes the matrix and stores the results in  $A$ . If  $m = n$  then the arguments MOVE,  $k$ , MDIM are ignored.

If  $m \neq n$  then  $k$  may be any integer. If  $k \leq 0$  then MOVE is ignored. Otherwise, if  $k \geq 1$  then MOVE is assumed to be an array of dimension  $k$ . MOVE is a storage area for saving information that may help speed up the transposition process. If no information is saved then TIP, DTIP, and CTIP may run 2-10 times slower than TPOSE, DTPOSE, and CTPOSE. However, the use of a storage area may or may not actually speed up the transposition process. This depends entirely on the values of  $m$  and  $n$ . MDIM is a variable that is set by the routine. After the routine terminates, MDIM will be the estimated optimum setting for  $k$  for the current values of  $m$  and  $n$ .

**Programming.** The routines TIP, DTIP, and CTIP employ the subroutine INFCTR. The routines were written by Norman Brenner (Dept. of Earth and Planetary Sciences, Massachusetts Institute of Technology), and modified by A. H. Morris.

**Reference.** Brenner, Norman, "Algorithm 467. Matrix Transposition in Place," *Comm. ACM* **16** (1973), pp. 692-694.

## COMPUTING ADJOINTS OF COMPLEX MATRICES

If  $A = (a_{ij})$  then let  $A^* = (\bar{a}_{ji})$  denote the adjoint of  $A$ . The following subroutines are available for computing  $A^*$ .

**CALL CMADJ**( $m, n, A, ka, B, kb$ )

$A$  is a complex  $m \times n$  matrix and  $B$  a complex 2-dimensional array. CMADJ computes  $A^*$  and stores the results in  $B$ . The input arguments  $ka$  and  $kb$  have the following values:

$ka$  = the number of rows in the dimension statement for  $A$  in the calling program

$kb$  = the number of rows in the dimension statement for  $B$  in the calling program

It is assumed that  $ka \geq m$  and  $kb \geq n$ , and that  $B$  contains at least  $n$  columns.

**Remark.** CMADJ combines the following operations:

CALL CTPOSE( $m, n, A, ka, B, kb$ )

CALL CMCONJ( $n, m, B, kb, B, kb$ )

It is assumed that  $A$  and  $B$  are different storage areas.

**Programmer.** A. H. Morris

**CALL CTRANS**( $ka, n, A$ )

$A$  is a complex array of dimension  $ka \times n$  which contains an  $n \times n$  matrix. CTRANS computes the adjoint of the matrix and stores the results in  $A$ . It is assumed that  $ka \geq n$ .

**Programmer.** George J. Davis (Georgia State University)

## MATRIX ADDITION

The matrix sum  $C = A + B$  can be computed by the following subroutines:

**CALL MADD( $m, n, A, ka, B, kb, C, kc$ )**

$A$  and  $B$  are real  $m \times n$  matrices and  $C$  a real 2-dimensional array. MADD computes  $A + B$  and stores the results in  $C$ . The input arguments  $ka, kb, kc$  have the following values:

$ka$  = the number of rows in the dimension statement for  $A$  in the calling program

$kb$  = the number of rows in the dimension statement for  $B$  in the calling program

$kc$  = the number of rows in the dimension statement for  $C$  in the calling program

It is assumed that  $ka \geq m, kb \geq m, kc \geq m$ , and that  $C$  contains at least  $n$  columns.

The array  $C$  may begin in the same location as  $A$  or  $B$ . If  $C$  begins in the same location as  $A$  then it is assumed that  $kc = ka$ . In this case, the result  $C$  will overwrite the input data  $A$ . Similarly, if  $C$  begins in the same location as  $B$  then it is assumed that  $kc = kb$ . Otherwise, if  $C$  does not begin in the same location as  $A$  or  $B$ , then it is assumed that the storage area for  $C$  does not overlap with the storage areas for  $A$  and  $B$ . In this case there is no restriction on  $kc$  (other than the customary restriction that  $kc \geq m$ ).

**Programmer.** A. H. Morris

**CALL SMADD( $n, A, B, C$ )**

$A$  and  $B$  are real packed  $n \times n$  symmetric matrices and  $C$  is a real array whose dimension is equal to or greater than  $n(n+1)/2$ . SMADD computes  $A + B$ , which is also a symmetric matrix, and stores the results in packed form in  $C$ .

The array  $C$  may begin in the same location as  $A$  or  $B$ . If  $C$  begins in the same location as  $A$  then the result  $C$  will overwrite the input data  $A$ . Similarly,  $B$  will be overwritten if  $C$  begins in the same location as  $B$ . Otherwise, if  $C$  does not begin in the same location as  $A$  or  $B$ , then it is assumed that the storage area for  $C$  does not overlap with the storage areas for  $A$  and  $B$ .

**Programmer.** A. H. Morris

**CALL DMADD( $m, n, A, ka, B, kb, C, kc$ )**

$A$  and  $B$  are double precision  $m \times n$  matrices and  $C$  a double precision 2-dimensional array. DMADD computes  $A + B$  and stores the results in  $C$ . The arguments  $ka, kb, kc$  have the following values:

$ka$  = the number of rows in the dimension statement for  $A$  in the calling program

$kb$  = the number of rows in the dimension statement for  $B$  in the calling program

$kc$  = the number of rows in the dimension statement for  $C$  in the calling program

It is assumed that  $ka \geq m, kb \geq m, kc \geq m$ , and that  $C$  contains at least  $n$  columns.

The array  $C$  may begin in the same location as  $A$  or  $B$ . If  $C$  begins in the same location as  $A$  then it is assumed that  $kc = ka$ . In this case, the result  $C$  will overwrite the input data  $A$ . Similarly, if  $C$  begins in the same location as  $B$  then it is assumed that  $kc = kb$ . Otherwise, if  $C$  does not begin in the same location as  $A$  or  $B$ , then it is assumed that the storage area for  $C$  does not overlap with the storage areas for  $A$  and  $B$ . In this case there is no restriction on  $kc$  (other than the customary restriction that  $kc \geq m$ ).

**Programmer.** A. H. Morris

**CALL CMADD**( $m, n, A, ka, B, kb, C, kc$ )

$A$  and  $B$  are complex  $m \times n$  matrices and  $C$  a complex 2-dimensional array. CMADD computes  $A + B$  and stores the results in  $C$ . The arguments  $ka, kb, kc$  have the following values:

$ka$  = the number of rows in the dimension statement for  $A$  in the calling program

$kb$  = the number of rows in the dimension statement for  $B$  in the calling program

$kc$  = the number of rows in the dimension statement for  $C$  in the calling program

It is assumed that  $ka \geq m, kb \geq m, kc \geq m$ , and that  $C$  contains at least  $n$  columns.

The array  $C$  may begin in the same location as  $A$  or  $B$ . If  $C$  begins in the same location as  $A$  then it is assumed that  $kc = ka$ . In this case, the result  $C$  will overwrite the input data  $A$ . Similarly, if  $C$  begins in the same location as  $B$  then it is assumed that  $kc = kb$ . Otherwise, if  $C$  does not begin in the same location as  $A$  or  $B$ , then it is assumed that the storage area for  $C$  does not overlap with the storage areas for  $A$  and  $B$ . In this case there is no restriction on  $kc$  (other than the customary restriction that  $kc \geq m$ ).

**Programmer.** A. H. Morris



## MATRIX SUBTRACTION

The matrix difference  $C = A - B$  can be computed by the following subroutines:

**CALL MSUBT**( $m, n, A, ka, B, kb, C, kc$ )

$A$  and  $B$  are real  $m \times n$  matrices and  $C$  a real 2-dimensional array. MSUBT computes  $A - B$  and stores the results in  $C$ . The input arguments  $ka, kb, kc$  have the following values:

$ka$  = the number of rows in the dimension statement for  $A$  in the calling program

$kb$  = the number of rows in the dimension statement for  $B$  in the calling program

$kc$  = the number of rows in the dimension statement for  $C$  in the calling program

It is assumed that  $ka \geq m, kb \geq m, kc \geq m$ , and that  $C$  contains at least  $n$  columns.

The array  $C$  may begin in the same location as  $A$  or  $B$ . If  $C$  begins in the same location as  $A$  then it is assumed that  $kc = ka$ . In this case, the result  $C$  will overwrite the input data  $A$ . Similarly, if  $C$  begins in the same location as  $B$  then it is assumed that  $kc = kb$ . Otherwise, if  $C$  does not begin in the same location as  $A$  or  $B$ , then it is assumed that the storage area for  $C$  does not overlap with the storage areas for  $A$  and  $B$ . In this case there is no restriction on  $kc$  (other than the customary restriction that  $kc \geq m$ ).

**Programmer.** A. H. Morris

**CALL SMSUBT**( $n, A, B, C$ )

$A$  and  $B$  are real packed  $n \times n$  symmetric matrices and  $C$  is a real array whose dimension is equal to or greater than  $n(n+1)/2$ . SMSUBT computes  $A - B$ , which is also a symmetric matrix, and stores the results in packed form in  $C$ .

The array  $C$  may begin in the same location as  $A$  or  $B$ . If  $C$  begins in the same location as  $A$  then the result  $C$  will overwrite the input data  $A$ . Similarly,  $B$  will be overwritten if  $C$  begins in the same location as  $B$ . Otherwise, if  $C$  does not begin in the same location as  $A$  or  $B$ , then it is assumed that the storage area for  $C$  does not overlap with the storage areas for  $A$  and  $B$ .

**Programmer.** A. H. Morris

**CALL DMSUBT**( $m, n, A, ka, B, kb, C, kc$ )

$A$  and  $B$  are double precision  $m \times n$  matrices and  $C$  a double precision 2-dimensional array. DMSUBT computes  $A - B$  and stores the results in  $C$ . The arguments  $ka, kb, kc$  have the following values:

$ka$  = the number of rows in the dimension statement for  $A$  in the calling program

$kb$  = the number of rows in the dimension statement for  $B$  in the calling program

$kc$  = the number of rows in the dimension statement for  $C$  in the calling program

It is assumed that  $ka \geq m, kb \geq m, kc \geq m$ , and that  $C$  contains at least  $n$  columns.

The array  $C$  may begin in the same location as  $A$  or  $B$ . If  $C$  begins in the same location as  $A$  then it is assumed that  $kc = ka$ . In this case, the result  $C$  will overwrite the input data  $A$ . Similarly, if  $C$  begins in the same location as  $B$  then it is assumed that  $kc = kb$ . Otherwise, if  $C$  does not begin in the same location as  $A$  or  $B$ , then it is assumed that the storage area for  $C$  does not overlap with the storage areas for  $A$  and  $B$ . In this case there is no restriction on  $kc$  (other than the customary restriction that  $kc \geq m$ ).

**Programmer.** A. H. Morris

**CALL CMSUBT**( $m, n, A, ka, B, kb, C, kc$ )

$A$  and  $B$  are complex  $m \times n$  matrices and  $C$  a complex 2-dimensional array. CMSUBT computes  $A - B$  and stores the results in  $C$ . The input arguments  $ka, kb, kc$  have the following values:

$ka$  = the number of rows in the dimension statement for  $A$  in the calling program

$kb$  = the number of rows in the dimension statement for  $B$  in the calling program

$kc$  = the number of rows in the dimension statement for  $C$  in the calling program

It is assumed that  $ka \geq m, kb \geq m, kc \geq m$ , and that  $C$  has at least  $n$  columns.

The array  $C$  may begin in the same location as  $A$  or  $B$ . If  $C$  begins in the same location as  $A$  then it is assumed that  $kc = ka$ . In this case, the result  $C$  will overwrite the input data  $A$ . Similarly, if  $C$  begins in the same location as  $B$  then it is assumed that  $kc = kb$ . Otherwise, if  $C$  does not begin in the same location as  $A$  or  $B$ , then it is assumed that the storage area for  $C$  does not overlap with the storage areas for  $A$  and  $B$ . In this case there is no restriction on  $kc$  (other than the customary restriction that  $kc \geq m$ ).

**Programmer.** A. H. Morris

## MATRIX MULTIPLICATION

The matrix product  $C = AB$  may be computed by the subroutines MTMS, DMTMS, CMTMS or MPROD, DMPROD, CMPROD. MTMS, DMTMS, and CMTMS can be used if the storage area for  $C$  does not overlap with the storage areas for  $A$  and  $B$ . Otherwise, if the components of  $C$  are to be stored in the storage area for  $A$  or  $B$ , then MPROD, DMPROD, or CMPROD must be used.

```
CALL MTMS( $m, n, \ell, A, ka, B, kb, C, kc$ )  
CALL DMTMS( $m, n, \ell, A, ka, B, kb, C, kc$ )  
CALL CMTMS( $m, n, \ell, A, ka, B, kb, C, kc$ )
```

MTMS is used if  $A$  and  $B$  are real matrices and  $C$  a real array, DMTMS is used if  $A$  and  $B$  are double precision matrices and  $C$  a double precision array, and CMTMS is used if  $A$  and  $B$  are complex matrices and  $C$  a complex array.

$A$  is a matrix having  $m$  rows and  $n$  columns,  $B$  a matrix having  $n$  rows and  $\ell$  columns, and  $C$  a 2-dimensional array. The routine computes the product  $AB$  and stores the results in  $C$ . The input arguments  $ka, kb, kc$  have the following values:

$ka$  = the number of rows in the dimension statement for  $A$  in the calling program

$kb$  = the number of rows in the dimension statement for  $B$  in the calling program

$kc$  = the number of rows in the dimension statement for  $C$  in the calling program

It is assumed that  $ka \geq m, kb \geq n, kc \geq m$ , and that  $C$  contains at least  $\ell$  columns.

**Remark.** It is assumed that the storage area for  $C$  is separate from the storage areas for  $A$  and  $B$ .

**Programmer.** A. H. Morris

```
CALL MPROD( $m, n, \ell, A, ka, B, kb, C, kc, WK$ )  
CALL DMPROD( $m, n, \ell, A, ka, B, kb, C, kc, WK$ )  
CALL CMPROD( $m, n, \ell, A, ka, B, kb, C, kc, WK$ )
```

MPROD is used if  $A$  and  $B$  are real matrices and  $C$  and  $WK$  are real arrays, DMPROD is used if  $A$  and  $B$  are double precision matrices and  $C$  and  $WK$  are double precision arrays, and CMPROD is used if  $A$  and  $B$  are complex matrices and  $C$  and  $WK$  are complex arrays.

$A$  is a matrix having  $m$  rows and  $n$  columns,  $B$  a matrix having  $n$  rows and  $\ell$  columns, and  $C$  a 2-dimensional array. The routine computes the product  $AB$  and stores the results in  $C$ . The input arguments  $ka, kb, kc$  have the following values:

$ka$  = the number of rows in the dimension statement for  $A$  in the calling program

$kb$  = the number of rows in the dimension statement for  $B$  in the calling program

$kc$  = the number of rows in the dimension statement for  $C$  in the calling program

It is assumed that  $ka \geq m, kb \geq n, kc \geq m$ , and that  $C$  contains at least  $\ell$  columns.

WK is an array that serves as a temporary storage area. The matrix  $C$  may begin in the same location as  $A$  or  $B$ . If  $C$  begins in the same location as  $A$ , then it is assumed that  $kc = ka$  and that the dimension of WK is equal to or greater than  $\ell$ . In this case, the result  $C$  will overwrite the input data  $A$ . Similarly, if  $C$  begins in the same location as  $B$  then it is assumed that  $kc = kb$  and that the dimension of WK is equal to or greater than  $m$ . Otherwise, if  $C$  does not begin in the same location as  $A$  or  $B$ , then it is assumed that the storage area for  $C$  does not overlap with the storage areas for  $A$  and  $B$ . In this case, the array WK is not referenced.

**Remark.** If  $C$  begins in the same location as  $A$  or  $B$ , then it is assumed that the storage areas for  $A$  and  $B$  are distinct storage areas.

**Programming.** MPROD employs the subroutines RLOC and YCHG, DMPROD employs the subroutines DLOC and DYCHG, and CMPROD employs the subroutines CLOC and CYCHG. The routines were written by A. H. Morris.

## PRODUCT OF A PACKED SYMMETRIC MATRIX AND A VECTOR

Let  $A$  denote a packed symmetric matrix of order  $n$  and  $x$  a column vector of dimension  $n$  where  $n \geq 1$ . Then the following subroutines are available for computing the product  $Ax$ .

```
CALL SVPRD( $A, n, x, y$ )  
CALL DSVPRD( $A, n, x, y$ )
```

The argument  $y$  is an array of dimension  $n$ . SVPRD is called when  $A, x, y$  are real arrays and DSVPRD is called when  $A, x, y$  are double precision arrays. When either of these routines is called,  $Ax$  is computed and stored in  $y$ .

Programmer. A. H. Morris

## TRANSPOSE MATRIX PRODUCTS

If  $A^t$  denotes the transpose of  $A$ , then the matrix product  $C = A^t B$  can be computed by the following subroutine:

**CALL TMPROD**( $m, n, \ell, A, ka, B, kb, C, kc$ )

$A$  is a real matrix having  $m$  rows and  $n$  columns,  $B$  a real matrix having  $m$  rows and  $\ell$  columns, and  $C$  a real 2-dimensional array. TMPROD computes  $A^t B$  and stores the results in  $C$ . The input arguments  $ka, kb, kc$  have the following values:

$ka$  = the number of rows in the dimension statement for  $A$  in the calling program

$kb$  = the number of rows in the dimension statement for  $B$  in the calling program

$kc$  = the number of rows in the dimension statement for  $C$  in the calling program

Here it is assumed that  $ka \geq m, kb \geq m, kc \geq n$ , and that  $C$  contains at least  $\ell$  columns. Also it is assumed that the storage area for  $C$  does not overlap with the storage areas for  $A$  and  $B$ .

**Note.** All inner products  $\sum_k a_{ki} b_{kj}$  are computed in double precision and the results stored in single precision.

**Programmer.** A. H. Morris

## SYMMETRIC MATRIX PRODUCTS

If  $A^t$  denotes the transpose of  $A$ , then the matrix product  $A^t A$  can be computed and its packed form stored in the array  $B$  by the following subroutine:

**CALL SMPROD**( $m, n, A, ka, B$ )

$A$  is a real  $m \times n$  matrix and  $B$  a real array whose dimension is equal to or greater than  $n(n+1)/2$ . SMPROD computes  $A^t A$  and stores its packed form in  $B$ . The input argument  $ka$  has the following value:

$ka$  = the number of rows in the dimension statement for  $A$  in the calling program  
It is assumed that  $ka \geq m$ .

**Note.** All inner products  $\sum_k a_{ki} a_{kj}$  are computed in double precision and the results stored in single precision.

**Programmer.** A. H. Morris

## KRONECKER PRODUCT OF MATRICES

If  $A$  is an  $m \times n$  matrix and  $B$  a  $k \times \ell$  matrix, then the *Kronecker product*  $A \otimes B$  is defined by

$$A \otimes B = \begin{pmatrix} a_{11}B & \cdots & a_{1n}B \\ \vdots & & \vdots \\ a_{m1}B & \cdots & a_{mn}B \end{pmatrix}.$$

From this definition we obtain:

- (1) (Transpose Equality)  $(A \otimes B)^t = A^t \otimes B^t$ .
- (2)  $(A \otimes B) \otimes E = A \otimes (B \otimes E)$  for any matrix  $E$ .
- (3)  $(A \otimes B)(C \otimes D) = (AC) \otimes (BD)$  if  $C$  is a matrix having  $n$  rows and  $D$  a matrix having  $\ell$  rows.

If  $A$  and  $B$  are complex square matrices of orders  $m$  and  $k$  respectively, then from the Jordan canonical forms of  $A$  and  $B$  the determinant equality  $\det(A \otimes B) = (\det A)^k (\det B)^m$  can be verified. Moreover, if  $A$  and  $B$  are nonsingular then  $(A \otimes B)^{-1} = A^{-1} \otimes B^{-1}$  from (3).

The following subroutines are available for computing  $C = A \otimes B$ .

```
CALL KPROD(A, ka, m, n, B, kb, k, l, C, kc)
CALL DKPROD(A, ka, m, n, B, kb, k, l, C, kc)
CALL CKPROD(A, ka, m, n, B, kb, k, l, C, kc)
```

It is assumed that  $A$  is an  $m \times n$  matrix,  $B$  a  $k \times \ell$  matrix, and  $C$  a 2-dimensional array. KPROD is used if  $A, B, C$  are real arrays, DKPROD is used if  $A, B, C$  are double precision arrays, and CKPROD is used if  $A, B, C$  are complex arrays. When the routine is called,  $A \otimes B$  is computed and stored in  $C$ .

The arguments  $ka$ ,  $kb$ , and  $kc$  have the following values:

$ka$  = the number of rows in the dimension statement for  $A$  in the calling program

$kb$  = the number of rows in the dimension statement for  $B$  in the calling program

$kc$  = the number of rows in the dimension statement for  $C$  in the calling program

It is assumed that  $ka \geq m$ ,  $kb \geq k$ ,  $kc \geq mk$ , and that  $C$  contains at least  $n\ell$  columns.

**Remark.** It is assumed that the array  $C$  does not overlap with  $A$  or  $B$ .

**Programmer.** A. H. Morris



## INVERTING GENERAL REAL MATRICES AND SOLVING GENERAL SYSTEMS OF REAL LINEAR EQUATIONS

The subroutines CROUT, KROUT, MSLV, NPIVOT, and DMSLV are available for inverting real matrices  $A$  and solving systems of real linear equations  $AX = B$ . CROUT, KROUT, MSLV, and NPIVOT solve single precision problems, and DMSLV solves double precision problems.

All the routines except NPIVOT are general-purpose, employing partial pivot Gauss elimination. NPIVOT can only occasionally be used since it uses Gauss-Jordan elimination with no pivot search. Normally, CROUT and KROUT produce the same results, which will be of equal or greater accuracy than the results produced by MSLV. However, since many of the calculations are performed in double precision in CROUT and KROUT, whereas only single precision is used in MSLV, MSLV may run 2-3 times faster than CROUT and KROUT.

**CALL CROUT(MO,  $n$ ,  $m$ ,  $A$ ,  $ka$ ,  $B$ ,  $kb$ ,  $D$ , INDEX, TEMP)**

$A$  is a real matrix of order  $n$  where  $n \geq 1$ . If  $MO = 0$  then the inverse of  $A$  is computed and stored in  $A$ . If  $MO \neq 0$  then the inverse is not computed.

The argument  $m$  is an integer. If  $m \geq 1$  then  $B$  is a real matrix having  $n$  rows and  $m$  columns. In this case the matrix equation  $AX = B$  is solved and the solution  $X$  is stored in  $B$ . If  $m \leq 0$  then there are no equations to be solved. In this case the argument  $B$  is ignored.

The argument  $ka$  is the number of rows in the dimension statement for  $A$  in the calling program, and  $kb$  the number of rows in the dimension statement for  $B$  in the calling program. If  $m \leq 0$  then the argument  $kb$  is ignored.

$D$  is a real variable. When CROUT is called,  $D$  is assigned the value  $\det(A)$  where  $\det(A)$  is the determinant of  $A$ . If  $D$  is found to have the value 0 then the routine immediately terminates.

INDEX is an array of dimension  $n - 1$  or larger that is used by the routine for keeping track of the row interchanges that are made. If  $MO \neq 0$  then INDEX is ignored.

TEMP is a real array of dimension  $n$  or larger that is a work space for the routine. If  $MO \neq 0$  then TEMP is ignored.

### Remarks.

- (1) KROUT should be used instead of CROUT when the determinant  $D$  is not needed. Underflow or overflow in the calculation of  $D$  may cause CROUT to terminate prematurely.
- (2) The matrix  $A$  is destroyed.

**Algorithm.** The Crout procedure is used. All inner products are computed in double precision and the results returned in single precision. Partial pivoting is performed.

**Programming.** CROUT calls the subroutine CROUT0. These routines were written by A. H. Morris.

**CALL KROUT**(MO,  $n$ ,  $m$ ,  $A$ ,  $ka$ ,  $B$ ,  $kb$ , IERR, INDEX, TEMP)

$A$  is a real matrix of order  $n$  where  $n \geq 1$ . If MO = 0 then the inverse of  $A$  is computed and stored in  $A$ . If MO  $\neq$  0 then the inverse is not computed.

The argument  $m$  is an integer. If  $m \geq 1$  then  $B$  is a real matrix having  $n$  rows and  $m$  columns. In this case the matrix equation  $AX = B$  is solved and the solution  $X$  is stored in  $B$ . If  $m \leq 0$  then there are no equations to be solved. In this case the argument  $B$  is ignored.

The argument  $ka$  is the number of rows in the dimension statement for  $A$  in the calling program, and  $kb$  the number of rows in the dimension statement for  $B$  in the calling program. If  $m \leq 0$  then the argument  $kb$  is ignored.

INDEX is an array of dimension  $n - 1$  or larger that is used by the routine for keeping track of the row interchanges that are made. If MO  $\neq$  0 then INDEX is ignored.

TEMP is a real array of dimension  $n$  or larger that is a work space for the routine. If MO  $\neq$  0 then TEMP is ignored.

**Error Return.** IERR is a variable that reports the status of the results. When the routine terminates IERR has one of the following values:

IERR = 0 The requested results were obtained.

IERR = -1 Either  $n$ ,  $ka$ , or  $kb$  is incorrect. In this case,  $A$  and  $B$  have not been modified.

IERR =  $k$  The  $k^{\text{th}}$  column of  $A$  has been reduced to a column containing only zeros.

When an error is detected, the routine immediately terminates.

**Remark.** The matrix  $A$  is destroyed.

**Algorithm.** The CROUT procedure is used. All inner products are computed in double precision and the results returned in single precision. Partial pivoting is performed.

**Programming.** KROUT calls the subroutine KROUT0. These routines were written by A. H. Morris.

**CALL NPIVOT**( $n$ ,  $m$ ,  $A$ ,  $ka$ ,  $B$ ,  $kb$ ,  $D$ , IERR)

$A$  is a real matrix of order  $n$  where  $n \geq 1$ . When NPIVOT is called the inverse of  $A$  is computed and stored in  $A$ .

The argument  $m$  is an integer. If  $m \geq 1$  then  $B$  is a real matrix having  $n$  rows and  $m$  columns. In this case the matrix equation  $AX = B$  is solved and the solution  $X$  is stored in  $B$ . If  $m \leq 0$  then there are no equations to be solved. In this case the argument  $B$  is ignored.

The argument  $ka$  is the number of rows in the dimension statement for  $A$  in the calling program, and  $kb$  the number of rows in the dimension statement for  $B$  in the calling program. If  $m \leq 0$  then  $kb$  is ignored.

$D$  is a real variable. On input  $D$  must be assigned a value by the user. If the input value is  $r$ , then when NPIVOT terminates  $D = rd$  where  $d$  is the determinant of  $A$ .

**Error Return.** IERR is an integer variable. If inversion is successful then IERR is assigned the value 0. Otherwise, IERR = 1 if NPIVOT cannot invert the matrix.

**Algorithm.** The Gauss-Jordan procedure is used. However, no pivot searching is done. NPIVOT terminates (with IERR set to 1) whenever a zero pivot element is encountered.

**Remarks.** Since pivot search is frequently needed to invert a matrix, and since pivot search is normally required to obtain accurate results, *NPIVOT should not be used except on those occasions when pivot search is known to be superfluous.*

**Programmer.** A. H. Morris

```
CALL MSLV(MO, n, m, A, ka, B, kb, D, RCOND, IERR, IPVT, WK)
CALL DMSLV(MO, n, m, A, ka, B, kb, D, RCOND, IERR, IPVT, WK)
```

$A$  is a matrix of order  $n$  where  $n \geq 1$ . If  $MO = 0$  then the inverse of  $A$  is computed and stored in  $A$ . If  $MO \neq 0$  then the inverse is not computed.

The argument  $m$  is an integer. If  $m \geq 1$  then  $B$  is a matrix having  $n$  rows and  $m$  columns. In this case the matrix equation  $AX = B$  is solved and the solution  $X$  is stored in  $B$ . If  $m \leq 0$  then there are no equations to be solved. In this case the argument  $B$  is ignored.

The argument  $ka$  is the number of rows in the dimension statement for  $A$  in the calling program, and  $kb$  the number of rows in the dimension statement for  $B$  in the calling program. If  $m \leq 0$  then the argument  $kb$  is ignored.

$D$  is an array of dimension 2. When either routine is called the determinant  $\det(A)$  of the matrix  $A$  is computed. If  $\det(A) = d \cdot 10^k$  where  $1 \leq |d| < 10$  and  $k$  an integer, then  $d$  is stored in  $D(1)$  and the exponent  $k$  is stored in floating point form in  $D(2)$ .

RCOND is a variable. When either routine is called, the routine makes an estimate  $c$  of the condition number of the matrix  $A$  (relative to the  $L_1$  norm). RCOND is assigned the value  $1/c$ .

IPVT is an integer array of dimension  $n$  or larger that is used by the routines for

keeping track of the row interchanges that are made. WK is an array of dimension  $n$  or larger that is used as a work space.

#### Remarks.

- (1) If MSLV is called then it is assumed that  $A$  and  $B$  are real matrices,  $D$  and WK real arrays, and RCOND a real variable. Otherwise, if DMSLV is called then it is assumed that  $A$  and  $B$  are double precision matrices,  $D$  and WK double precision arrays, and RCOND a double precision variable.
- (2) RCOND satisfies  $0 \leq \text{RCOND} \leq 1$ . If  $\text{RCOND} \approx 10^{-k}$  then one can expect the results to have approximately  $k$  fewer significant digits of accuracy than the elements of  $A$ . For example, if MSLV is used to invert a matrix in the 14 digit CDC single precision arithmetic and  $\text{RCOND} = .4\text{E}-3$ , then the computed coefficients of the inverse matrix should normally be accurate to about 11 digits. In general, RCOND characterizes how well or poorly conditioned the problem is. If  $\text{RCOND} \approx 1$  then one should expect the results to be almost as accurate as the original data  $A$ . However, if  $\text{RCOND} \approx 0$  then one should expect the results to be nonsense.
- (3) The matrix  $A$  is destroyed.

**Error Return.** IERR is an integer variable. If RCOND is sufficiently large so that  $1 + \text{RCOND} > 1$ , then IERR is set to 0 and the problem is solved. Otherwise, if  $1 + \text{RCOND} = 1$  then IERR is set to 1 and the routine terminates. In this case,  $A$  will have been destroyed but  $B$  will not have been modified. Also the determinant will not have been computed.

**Algorithm.** The partial pivot Gauss elimination procedure is used.

**Programming.** MSLV and DMSLV are driver routines for the LINPACK subroutines SGECO, SGEFA, SGESL, SGEDI and DGECO, DGEFA, DGESL, DGEDI. The subroutines were coded by Cleve Moler (University of New Mexico). The subroutines employ the vector routines SSWAP, SDOT, SSCAL, SAXPY, SASUM, ISAMAX and DSWAP, DDOT, DSCAL, DAXPY, DASUM, IDAMAX.

#### References.

- (1) Dongarra, J. J., Bunch, J. R., Moler, C. B., and Stewart, G. W., *LINPACK Users' Guide*, Society for Industrial and Applied Mathematics, Philadelphia, 1979.
- (2) Cline, A. K., Moler, C. B., Stewart, G. W., and Wilkinson, J. H., "An Estimate for the Condition Number of a Matrix," *SIAM Journal of Numerical Analysis* **16** (1979), pp. 368-375.

## SOLUTION OF REAL EQUATIONS WITH ITERATIVE IMPROVEMENT

Given a real  $n \times n$  matrix  $A$  and column vector  $b$ . The following subroutine is available for solving the equation  $Ax = b$ . Iterative improvement is performed to compute the solution  $x$  to machine accuracy.

**CALL SLVMP(MO,n, A, ka, b, X,WK,IWK,IND)**

MO is an input argument which specifies if SLVMP is being called for the first time. On an initial call, MO = 0 and we have the following setup:

A is a 2-dimensional real array of dimension  $ka \times n$  containing the matrix  $A$ ,  $b$  a real vector of dimension  $n$ , and  $X$  a real array of dimension  $n$ . When SLVMP is called,  $Ax = b$  is solved and the solution stored in  $X$ .  $A$  and  $b$  are not modified by the routine.

WK is a real array of dimension  $n^2 + n$  or larger, and IWK an integer array of dimension  $n$  or larger. These arrays are for internal use by the routine. On an initial call to SLVMP, an LU decomposition is obtained for  $A$  and stored in WK and IWK. Then the equation  $Ax = b$  is solved.

IND is an integer variable that reports the status of the results. On an initial call to SLVMP, when the routine terminates IND has one of the following values:

IND = 0 The solution  $X$  was obtained to machine accuracy.

IND = 1  $X$  was obtained, but not to machine accuracy.

IND =  $-k$  The  $k^{th}$  column of  $A$  was reduced to a column containing only zeros.

When IND =  $-k$ , no solution is obtained.

After an initial call to SLVMP, if IND = 0 or 1 on output, then the routine may be called to solve a new set of equations  $Ax = \tilde{b}$  without having to redecompose the matrix  $A$ . In this case, the input argument MO may be set to any nonzero value. When MO  $\neq$  0 it is assumed that only  $b$  has been modified. The routine employs the LU decomposition obtained on the initial call to SLVMP to solve the new set of equations  $Ax = b$ . On output  $X$  will contain the solution to the new set of equations. As before,  $A$  and  $b$  are not modified by the routine.

If SLVMP is recalled with MO  $\neq$  0, then when the routine terminates IND has one of the following values:

IND = 0 The solution  $X$  was obtained to machine accuracy.

IND = 1  $X$  was obtained, but not to machine accuracy.

**Programming.** SLVMP calls the subroutine LUIMP. These routines were written by A. H. Morris. The subroutines MCOPY, SGEFA, SGESL, SSCAL, SAXPY and functions SPM-PAR, SDOT, ISAMAX are also employed.

## SOLUTION OF ALMOST BLOCK DIAGONAL SYSTEMS OF LINEAR EQUATIONS

Consider a system  $Ax = b$  of linear equations where  $A$  is an  $n \times n$  matrix having the block structure

$$A = \begin{pmatrix} \boxed{A_1} & & & & 0 \\ & \boxed{A_2} & & & \\ & & \boxed{A_3} & & \\ & & & \ddots & \\ 0 & & & & \boxed{A_m} \end{pmatrix}$$

Here it is assumed that  $A_i$  is an  $r_i \times c_i$  matrix for  $i = 1, \dots, m$ , and that  $A_i$  and  $A_{i+1}$  may have  $\delta_i \geq 0$  columns in common for  $i < m$ . Thus  $\sum_{i=1}^m r_i = n$  and block  $A_i$  begins in column  $\sum_{k=1}^{i-1} (c_k - \delta_k) + 1$  for  $i \geq 2$ . It is also assumed that three successive blocks  $A_{i-1}, A_i, A_{i+1}$  do not have columns in common. Thus  $\delta_{i-1} + \delta_i \leq c_i$  for  $i = 2, \dots, m-1$ . If  $m \geq 2$  then the following subroutines are available for solving  $Ax = b$ .

**CALL ARCECO( $n, S, MTR, m, IWK, b, X, IND$ )**

$S$  is an array of dimension  $\sum_{i=1}^m r_i c_i$  or larger. On input  $S$  contains the blocks  $A_1, \dots, A_m$  of the matrix  $A$ . The blocks are stored in sequence.  $A_1$  is stored in the first  $r_1 c_1$  locations of  $S$ ,  $A_2$  is stored in the next  $r_2 c_2$  locations, etc. For each  $A_i$ , the columns of  $A_i$  are stored in sequence in  $S$ .

**Example.**

$$\text{If } A = \begin{pmatrix} \boxed{a_{11} \ a_{12} \ a_{13}} & 0 & 0 \\ \boxed{a_{21} \ a_{22} \ a_{23}} & 0 & 0 \\ 0 & \boxed{a_{32} \ a_{33} \ a_{34}} & 0 \\ 0 & \boxed{a_{42} \ a_{43} \ a_{44}} & 0 \\ 0 & 0 & 0 & 0 & \boxed{a_{55}} \end{pmatrix}$$

then  $n = 5$ ,  $m = 3$ ,  $\delta_1 = 2$ , and  $\delta_2 = 0$ . Also,  $S$  is an array containing the data  $a_{11}, a_{21}, a_{12}, a_{22}, a_{13}, a_{23}, a_{32}, a_{42}, a_{33}, a_{43}, a_{34}, a_{44}, a_{55}$ .

$MTR$  is an integer matrix of dimension  $3 \times m$  containing the block information of the matrix  $A$ :

$$MTR(1, i) = r_i \quad (i = 1, \dots, m)$$

$$MTR(2, i) = c_i \quad (i = 1, \dots, m)$$

$$MTR(3, i) = \delta_i \quad (i = 1, \dots, m-1)$$

For convenience, the routine sets  $MTR(3, m) = 0$ .

$X$  is an array of dimension  $n$  or larger. When ARCECO is called,  $A$  is decomposed and then the equations  $Ax = b$  are solved. The decomposition of  $A$  is stored in  $S$ , overwriting the initial input data  $A$ , and the solution  $x$  is stored in  $X$ . The vector  $b$  is destroyed by the routine.

IWK is an array of dimension  $n$  or larger for internal use by the routine. The pivot indices involved in the decomposition of  $A$  are stored in IWK.

IND is a variable that reports the status of the results. When ARCECO terminates, IND has one of the following values:

IND = 0 The system of equations was solved.

IND = 1 (Input Error) Either  $n$ ,  $m$ , or MTR is incorrect, or three successive blocks  $A_{i-1}$ ,  $A_i$ ,  $A_{i+1}$  of  $A$  have columns in common.

IND = -1  $A$  is a singular matrix. The equations cannot be solved.

**Usage.** After a call to ARCECO, if IND = 0 on output then the subroutine ARCESL (see below) may be called to solve a new set of equations  $Ax = \tilde{b}$  without having to redecompose the matrix  $A$ . ARCESL employs the decomposition of  $A$  obtained by ARCECO.

**Algorithm.** A modification of the alternate row and column elimination procedure by Varah is used.

**Programming.** ARCECO employs the routines ARCEDC, ARCEPR, ARCEPC, ARCESL, ARCEFS, ARCEFM, ARCEFE, ARCEBS, ARCEBM, and ARCEBE. These routines were developed by J. C. Diaz (Mobil Research and Development Corp., Farmers Branch, Texas), G. Fairweather (University of Kentucky), and P. Keast (University of Toronto).

**Reference.** Diaz, J. C., Fairweather, G., and Keast, P., "FORTRAN Packages for Solving Certain Almost Block Diagonal Linear Systems by Modified Alternate Row and Column Elimination," *ACM Trans. Math Software* 9 (1983), pp. 358-375.

#### CALL ARCESL( $S$ ,MTR, $m$ ,IWK, $b$ , $X$ )

The argument  $m$  is the number of blocks  $A_1, \dots, A_m$  in the matrix  $A$ . ARCESL may be used only when IND = 0 on output from ARCECO. In this case,  $S$  contains the decomposition of  $A$  obtained by ARCECO and IWK contains the pivot indices of the decomposition. The argument  $b$  is a vector of dimension  $n$ , and  $X$  an array of dimension  $n$  or larger. When ARCESL is called, the equations  $Ax = b$  are solved and the solution stored in  $X$ . The vector  $b$  is destroyed by the routine.

**Programming.** ARCESL calls the subroutines ARCEFS, ARCEFM, ARCEFE, ARCEBS, ARCEBM, and ARCEBE. These routines were developed by J. C. Diaz (Mobil Research and Development Corp., Farmers Branch, Texas), G. Fairweather (University of Kentucky), and P. Keast (University of Toronto).

## SOLUTION OF ALMOST BLOCK TRIDIAGONAL SYSTEMS OF LINEAR EQUATIONS

Consider a system  $Tx = b$  of linear equations where  $T$  is a square matrix having the block structure

$$T = \begin{pmatrix} A_1 & B_1 & C_1 & & \\ C_2 & A_2 & B_2 & & 0 \\ & C_3 & A_3 & B_3 & \\ 0 & \ddots & \ddots & \ddots & \\ & & C_{n-1} & A_{n-1} & B_{n-1} \\ & & B_n & C_n & A_n \end{pmatrix}.$$

Here it is assumed that  $A_k, B_k, C_k$  ( $k = 1, \dots, n$ ) are  $m \times m$  matrices, and that  $b$  is a column vector of dimension  $mn$ . If  $n \geq 4$  then the following subroutine is available for solving  $Tx = b$ .

**CALL BTSLV(MO,m,n,A,B,C,X,IP,IND)**

MO is an input argument which specifies if BTSLV is being called for the first time. On an initial call, MO = 0 and we have the following setup:

$A, B, C$  are 3-dimensional arrays of dimension  $m \times m \times n$  where the  $(i, j)$ -th elements of the matrices  $A_k, B_k, C_k$  are stored in  $A(i, j, k)$ ,  $B(i, j, k)$ ,  $C(i, j, k)$  for  $k = 1, \dots, n$ .  $A, B, C$  are modified by the routine.

$X$  is an array of dimension  $mn$  or larger. On input, the vector  $b$  is stored in  $X$ . When BTSLV is called, if a solution  $x$  is obtained for  $Tx = b$  then the solution is stored in  $X$ .

IP is an array of dimension  $mn$  or larger that is used by the routine for listing the row interchanges that are made.

On an initial call to the routine, a block LU decomposition is performed on  $T$ , the results of which are stored in  $A, B, C$ . This decomposition involves row interchanges only within the diagonal blocks  $A_k$ ; i.e., no row interchanges are performed between rows of different blocks  $A_k$  and  $A_l$ . Thus it may occur that the decomposition of a nonsingular matrix  $T$  cannot be completed. IND is a variable that reports the status of the results. When BTSLV terminates, IND has one of the following values:

IND = 0  $T$  was decomposed and the equations  $Tx = b$  solved.

IND = -1 (Input Error) Either  $m < 1$  or  $n < 4$ .

IND =  $k$  The decomposition process failed in the  $k^{th}$  diagonal block. The routine cannot solve the equations in their present form.

After an initial call to BTSLV, if IND = 0 then the routine may be recalled with MO  $\neq 0$  and a new  $b$ . When MO  $\neq 0$ , then it is assumed that  $A, B, C$ , and IP have not been modified and that  $X$  contains the new  $b$ . The routine retrieves from  $A, B, C$ , and IP the block decomposition that was obtained on the initial call to BTSLV, and solves the new



system of equations  $Tx = b$ . The solution is stored in  $X$ . In this case, IND is not referenced by the routine.

**Programming.** BTSLV employs the subroutines DECBT, SOLBT, DEC, and SOL. These subroutines were written by Alan C. Hindmarsh (Lawrence Livermore Laboratory).

**Reference.** Hindmarsh, A. C., *Solution of Block-Tridiagonal Systems of Linear Algebraic Equations*, Report UCID-30150, Lawrence Livermore Laboratory, 1977.

## INVERTING SYMMETRIC REAL MATRICES AND SOLVING SYMMETRIC SYSTEMS OF REAL LINEAR EQUATIONS

The subroutines SMSLV and DSMSLV are available for inverting symmetric real matrices  $A$  and solving systems of real linear equations  $AX = B$ . SMSLV handles single precision problems and DSMSLV handles double precision problems. It is assumed that the matrix  $A$  is in packed form. If the inverse of  $A$  is computed, then the inverse is a symmetric matrix which will be stored in packed form.

**Note.** All eigenvalues of a real symmetric matrix  $A$  are real. The *inertia* of  $A$  is the ordered triple  $(\pi, \nu, \zeta)$  where  $\pi$  is the number of positive eigenvalues of  $A$ ,  $\nu$  the number of negative eigenvalues of  $A$ , and  $\zeta$  the number of zero eigenvalues of  $A$ . Thus, if  $n$  is the order of  $A$  then  $\pi + \nu + \zeta = n$ . Also  $A$  is positive definite (positive semi-definite, negative definite, negative semi-definite) if  $\pi = n$  ( $\nu = 0, \nu = n, \pi = 0$ ).

```
CALL SMSLV(MO,n,m,A,B,kb,D,RCOND,INERT,IERR,IPVT,WK)
CALL DSMSLV(MO,n,m,A,B,kb,D,RCOND,INERT,IERR,IPVT,WK)
```

$A$  is an array of dimension  $n(n+1)/2$  containing the elements of a packed  $n \times n$  symmetric matrix where  $n \geq 1$ . If  $MO = 0$  then the inverse of  $A$  is computed and stored in  $A$  in packed form. If  $MO \neq 0$  then the inverse of  $A$  is not computed.

The argument  $m$  is an integer. If  $m \geq 1$  then  $B$  is a matrix having  $n$  rows and  $m$  columns. In this case  $AX = B$  is solved and the solution  $X$  is stored in  $B$ . If  $m \leq 0$  then there are no equations to be solved. In this case  $B$  is ignored.

The argument  $kb$  is the number of rows in the dimension statement for  $B$  in the calling program. If  $m \leq 0$  then  $kb$  is ignored.

$D$  is an array of dimension 2. When either routine is called the determinant  $\det(A)$  of the matrix  $A$  is computed. If  $\det(A) = d \cdot 10^k$  where  $1 \leq |d| < 10$  and  $k$  an integer, then  $d$  is stored in  $D(1)$  and the exponent  $k$  is stored in floating point form in  $D(2)$ .

RCOND is a variable. When either routine is called, the routine makes an estimate  $c$  of the condition number of the matrix  $A$  (relative to the  $L_1$  norm). RCOND is assigned the value  $1/c$ .

INERT is an integer array of dimension 3. When either routine is called the inertia of the matrix  $A$  is computed. INERT(1) is set to the number of positive eigenvalues of  $A$ , INERT(2) is set to the number of negative eigenvalues, and INERT(3) is set to the number of zero eigenvalues.

IPVT is an integer array of dimension  $n$  or larger that is used by the routines for keeping track of the row and column interchanges that are made. WK is array of dimension  $n$  or larger that is used as a work space.

#### Remarks.

- (1) If SMSLV is called then it is assumed that  $A$  and  $B$  are real matrices,  $D$  and  $WK$  are real arrays, and  $RCOND$  is a real variable. Otherwise, if DSMSLV is called then it is assumed that  $A$  and  $B$  are double precision matrices,  $D$  and  $WK$  are double precision arrays, and  $RCOND$  is a double precision variable.
- (2)  $RCOND$  satisfies  $0 \leq RCOND \leq 1$ . If  $RCOND \approx 10^{-k}$  then one can expect the results to have approximately  $k$  fewer significant digits of accuracy than the elements of  $A$ . For example, if SMSLV is used to invert a matrix in the 14 digit CDC single precision arithmetic and  $RCOND = .4E-3$ , then the computed coefficients of the inverse matrix should normally be accurate to about 11 digits. In general,  $RCOND$  characterizes how well or poorly conditioned the problem is. If  $RCOND \approx 1$  then one should expect the results to be almost as accurate as the original data  $A$ . However, if  $RCOND \approx 0$  then one should expect the results to be nonsense.
- (3) The data in  $A$  is destroyed.

**Error Return.** IERR is an integer variable. If  $RCOND$  is sufficiently large so that  $1 + RCOND > 1$ , then IERR is set to 0 and the problem is solved. Otherwise, if  $1 + RCOND = 1$  then IERR is set to 1 and the routine terminates. In this case,  $A$  will have been destroyed but  $B$  will not have been modified. Also the determinant and inertia will not have been computed.

**Algorithm.** The diagonal pivoting factorization procedure is used. Partial pivoting is employed.

**Precision.** SMSLV and the more general routine MSLV have approximately the same accuracy, and DSMSLV and DMSLV have approximately the same accuracy.

**Efficiency.** Even though SMSLV performs approximately half the number of multiplications and divisions as MSLV, normally one can expect SMSLV to take about 70-80% of the time required by MSLV. However, for sparse matrices SMSLV may be 20-30% slower than MSLV. Similar comments hold for DSMSLV and DMSLV.

**Programming.** SMSLV and DSMSLV are driver routines for the LINPACK subroutines SSPCO, SSPFA, SSPSL, SSPDI and DSPCO, DSPFA, DSPSL, DSPDI. SSPCO and DSPCO were written by Cleve Moler (University of New Mexico). The remaining LINPACK subroutines were written by James Bunch (University of California, San Diego). The subroutines employ the vector routines SCOPY, SSWAP, SDOT, SSCAL, SAXPY, SASUM, ISAMAX and DSWAP, DDOT, DSCAL, DAXPY, DASUM, IDAMAX.

#### References.

- (1) Bunch, J. R. and Parlett, B. N., "Direct Methods for Solving Symmetric Indefinite Systems of Linear Equations," *SIAM J. Numerical Analysis* 8 (1971), pp. 639-655.
- (2) Bunch, J. R., "Analysis of the Diagonal Pivoting Method," *SIAM J. Numerical Analysis* 8 (1971), pp. 656-680.

- (3) Bunch, J. R., Kaufman, L., and Parlett, B.N., "Decomposition of a Symmetric Matrix," *Numerische Mathematik* **27** (1976), pp. 95-109.
- (4) Bunch, J., and Kaufman, L., "Some Stable Methods for Calculating Inertia and Solving Symmetric Linear Systems," *Math. Comp.* **31** (1977), pp. 163-179.
- (5) Cline, A. K., Moler, C. B., Stewart, G. W., and Wilkinson, J. H., "An Estimate for the Condition Number of a Matrix," *SIAM Numerical Analysis* **16** (1979), pp. 368-375.
- (6) Dongarra, J. J., Bunch, J. R., Moler, C. B., and Stewart, G. W., *LINPACK Users' Guide*, Society for Industrial and Applied Mathematics, Philadelphia, 1979.

## INVERTING POSITIVE DEFINITE SYMMETRIC MATRICES AND SOLVING POSITIVE DEFINITE SYMMETRIC SYSTEMS OF LINEAR EQUATIONS

The subroutines PCHOL and DPCHOL are available for inverting positive definite symmetric real matrices  $A$  and solving systems of real linear equations  $AX = B$ . PCHOL handles single precision problems and DPCHOL handles double precision problems. It is assumed that the matrix  $A$  is in packed form. If the inverse of  $A$  is computed then the inverse is a symmetric matrix which will be stored in packed form.

CALL PCHOL(MO,n,m,A,B,kb,IERR)  
CALL DPCHOL(MO,n,m,A,B,kb,IERR)

$A$  is an array of dimension  $n(n+1)/2$  or larger containing the elements of a packed  $n \times n$  positive definite symmetric matrix where  $n \geq 1$ . If  $MO = 0$  then the inverse of  $A$  is computed and stored in  $A$  in packed form. If  $MO \neq 0$  then the inverse is not computed.

The argument  $m$  is an integer. If  $m \geq 1$  then  $B$  is a matrix having  $n$  rows and  $m$  columns. In this case  $AX = B$  is solved and the solution  $X$  is stored in  $B$ . If  $m \leq 0$  then there are no equations to be solved. In this case  $B$  is ignored.

The argument  $kb$  is the number of rows in the dimension statement for  $B$  in the calling program. If  $m \leq 0$  then  $kb$  is ignored.

### Remarks.

- (1) If PCHOL is called then it is assumed that  $A$  and  $B$  are real arrays, and if DPCHOL is called then it is assumed that  $A$  and  $B$  are double precision arrays.
- (2) The data in  $A$  is destroyed.

**Error Return.** IERR is an integer variable. If  $A$  is positive definite then IERR is set to 0 and the problem is solved. Otherwise,  $IERR = k$  if the leading  $k \times k$  submatrix of  $A$  is not positive definite.

**Algorithm.** The Cholesky procedure is used.

**Precision.** The results obtained by PCHOL and DPCHOL are occasionally less accurate (up to 1 significant digit) than the results obtained by SMSLV and DSMSLV.

**Programming.** PCHOL and DPCHOL are driver routines for the LINPACK subroutines SPPFA, SPPSL, SPPDI and DPPFA, DPPSL, DPPDI. These subroutines were written by Cleve Moler (University of New Mexico). The functions SDOT, DDOT and subroutines SAXPY, SSCAL, DAXPY, DSCAL are also used.

**Reference.** Dongarra, J. J., Bunch, J. R., Moler, C. B., and Stewart, G. W., *LINPACK Users' Guide*, Society for Industrial and Applied Mathematics, Philadelphia, 1979.

## SOLUTION OF TOEPLITZ SYSTEMS OF LINEAR EQUATIONS

An  $n \times n$  Toeplitz matrix is a matrix of the form

$$\begin{pmatrix} a_0 & a_{-1} & a_{-2} & \dots & a_{-n+1} \\ a_1 & a_0 & a_{-1} & \dots & a_{-n+2} \\ a_2 & a_1 & a_0 & \dots & a_{-n+3} \\ \vdots & \vdots & \vdots & & \vdots \\ a_{n-1} & a_{n-2} & a_{n-3} & \dots & a_0 \end{pmatrix}.$$

For convenience, we denote this matrix by  $A_n$  for  $n \geq 1$ . If  $A_n$  is a real matrix where  $a_0 \neq 0$  and  $b$  is a real column vector, then the subroutines TOPLX and DTOPLX are available for solving the system of equations  $A_n x = b$ . TOPLX is a single precision routine and DTOPLX a double precision routine.

```
CALL TOPLX(A,b,x,n,G,H,IERR)
CALL DTOPLX(A,b,x,n,G,H,IERR)
```

TOPLX is used if  $A, b, x, G, H$  are real arrays, and DTOPLX is used if  $A, b, x, G, H$  are double precision arrays.

$A$  is an array of dimension  $2n-1$  or larger containing the coefficients  $A(j) = a_{j-n}$  ( $j = 1, \dots, 2n-1$ ) of the matrix  $A_n$ . The argument  $x$  is an array of dimension  $n$  or larger, and IERR is an integer variable. When the routine is called, if  $A_n x = b$  is solved then IERR is set to 0 and the solution is stored in  $x$ .  $A$  and  $b$  are not modified by the routine.

$G$  and  $H$  are arrays of dimension  $n$  or larger that are work spaces for the routine.

**Error Return.** IERR = 1 if the equations cannot be solved by the routine.

**Remarks.** The Levinson bordering procedure is used. This procedure is inductive, beginning with the solution  $x_1 = b_1/a_0$  of the equation  $a_0 x_1 = b_1$ . Given a solution for  $\sum_{j=1}^N a_{i-j} x_j = b_i$  ( $i = 1, \dots, N$ ) where  $N < n$ , then a solution for  $\sum_{j=1}^{N+1} a_{i-j} \bar{x}_j = b_i$  ( $i = 1, \dots, N+1$ ) is obtained where  $\bar{x}_1, \dots, \bar{x}_{N+1}$  is computed from  $x_1, \dots, x_N$ . This procedure fails when some submatrix  $A_N$  is singular (e.g., when  $a_0 = 0$ ). Also, the procedure may yield poor results when some  $A_N$  is exceedingly ill-conditioned. In such situations one must use a more general equation solver. In TOPLX and DTOPLX  $4(n-1)^2$  floating additions and  $4n(n-1)$  integer/floating multiplications and divisions are used when  $n \geq 2$ . Consequently, these routines are considerably more efficient than general equation solvers such as KROUT and DMSLV, but more restrictive and frequently less accurate.

**Programming.** TOPLX and DTOPLX are modifications by A. H. Morris of the subroutine TOEPLZ, written by George Rybicki.

**Reference.** Press, W. H., Flannery, B. P., Teukolsky, S. A., and Vetterling, W. T., *Numerical Recipes: The Art of Scientific Computing*, Cambridge University Press, 1986, pp. 47-52.

## INVERTING GENERAL COMPLEX MATRICES AND SOLVING GENERAL SYSTEMS OF COMPLEX LINEAR EQUATIONS

The subroutines CMSLV, CMSLV1, and DCMSLV are available for inverting complex matrices and solving systems of complex linear equations. CMSLV and CMSLV1 solve single precision problems and DCMSLV solves double precision problems.

**CALL CMSLV(MO,n,m,A,ka,B,kb,D,RCOND,IERR,IPVT,WK)**

$A$  is a complex matrix of order  $n$  where  $n \geq 1$ . If  $MO = 0$  then the inverse of  $A$  is computed and stored in  $A$ . If  $MO \neq 0$  then the inverse is not computed.

The argument  $m$  is an integer. If  $m \geq 1$  then  $B$  is a complex matrix having  $n$  rows and  $m$  columns. In this case the matrix equation  $AX = B$  is solved and the solution  $X$  is stored in  $B$ . If  $m \leq 0$  then there are no equations to be solved. In this case the argument  $B$  is ignored.

The argument  $ka$  is the number of rows in the dimension statement for  $A$  in the calling program, and the argument  $kb$  is the number of rows in the dimension statement for  $B$  in the calling program. If  $m \leq 0$  then the argument  $kb$  is ignored.

$D$  is a complex array of dimension 2. When CMSLV is called the determinant  $\det(A)$  of the matrix  $A$  is computed. If  $\det(A) = d \cdot 10^k$  where  $1 \leq |\operatorname{Re}(d)| + |\operatorname{Im}(d)| < 10$  and  $k$  an integer, then  $d$  is stored in  $D(1)$  and the exponent  $k$  is stored as a complex number in  $D(2)$ .

RCOND is a real variable. When CMSLV is called, the routine makes an estimate  $c$  of the condition number of the matrix  $A$  (relative to the modified  $L_1$  norm where each absolute value  $|z|$  is replaced with  $|\operatorname{Re}(z)| + |\operatorname{Im}(z)|$ ). RCOND is assigned the value  $1/c$ .

IPVT is an integer array of dimension  $n$  or larger that is used by the routine for keeping track of the row interchanges that are made. WK is a complex array of dimension  $n$  or larger that is used as a work space.

### Remarks.

- (1) RCOND satisfies  $0 \leq \text{RCOND} \leq 1$ . If  $\text{RCOND} \approx 10^{-k}$  then one can expect the results to have approximately  $k$  fewer significant digits of accuracy than the elements of  $A$ . For example, if CMSLV is used to invert a matrix in the 14 digit CDC single precision arithmetic and  $\text{RCOND} = .4E-3$ , then the computed coefficients of the inverse matrix should normally be accurate to about 11 digits. In general, RCOND characterizes how well or poorly conditioned the problem is. If  $\text{RCOND} \approx 1$  then one should expect the results to be almost as accurate as the original data  $A$ . However, if  $\text{RCOND} \approx 0$  then one should expect the results to be nonsense.
- (2) The matrix  $A$  is always destroyed.

**Error Return.** IERR is an integer variable. If RCOND is sufficiently large so that  $1 +$

RCOND > 1, then IERR is set to 0 and the problem is solved. Otherwise, if  $1 + \text{RCOND} = 1$  then IERR is set to 1 and the routine terminates. In this case,  $A$  will have been destroyed but  $B$  will not have been modified. Also the determinant will not have been computed.

**Algorithm.** The partial pivot Gauss elimination procedure is used. The pivots  $a_{kj}$  are selected so that  $|\text{Re}(a_{kj})| + |\text{Im}(a_{kj})| = \max \{ |\text{Re}(a_{ij})| + |\text{Im}(a_{ij})| : i = j, \dots, n \}$ .

**Programming.** CMSLV calls the LINPACK subroutines CGECO, CGEFA, CGESL, and CGEDI. These subroutines were written by Cleve Moler (University of New Mexico). The subroutines CSWAP, CSCAL, CSSCAL, CAXPY and functions CDOTC, SCASUM, ICA-MAX are also used.

# References.

- (1) Dongarra, J. J., Bunch, J. R., Moler, C. B., and Stewart, G. W., *LINPACK Users' Guide*, Society for Industrial and Applied Mathematics, Philadelphia, 1979.
- (2) Cline, A. K., Moler, C. B., Stewart, G. W., and Wilkinson, J. H., "An Estimate for the Condition Number of a Matrix," *SIAM Journal of Numerical Analysis* 16 (1979), pp. 368-375.

**CALL CMSLV1(MO,  $n$ ,  $m$ ,  $A$ ,  $ka$ ,  $B$ ,  $kb$ , IERR, IPVT, WK)**

$A$  is a complex matrix of order  $n$  where  $n \geq 1$ . If  $\text{MO} = 0$  then the inverse of  $A$  is computed and stored in  $A$ . If  $\text{MO} \neq 0$  then the inverse is not computed.

The argument  $m$  is an integer. If  $m \geq 1$  then  $B$  is a complex matrix having  $n$  rows and  $m$  columns. In this case the matrix equation  $AX = B$  is solved and the solution  $X$  is stored in  $B$ . If  $m \leq 0$  then there are no equations to be solved. In this case the argument  $B$  is ignored.

The argument  $ka$  is the number of rows in the dimension statement for  $A$  in the calling program, and the argument  $kb$  is the number of rows in the dimension statement for  $B$  in the calling program. If  $m \leq 0$  then the argument  $kb$  is ignored.

IPVT is an integer array of dimension  $n$  or larger that is used by the routines for keeping track of the row interchanges that are made.

WK is a complex array of dimension  $n$  or larger that is a work space for the routine. If  $\text{MO} \neq 0$  then WK is ignored.

**Error Return.** IERR is a variable that reports the status of the results. When CMSLV1 terminates IERR has one of the following values:

- IERR = 0 The requested results were obtained.
- IERR = -1 Either  $n$ ,  $ka$ , or  $kb$  is incorrect.
- IERR =  $k$  The  $k^{\text{th}}$  column of  $A$  has been reduced to a column containing zeros. The requested results cannot be obtained.



**Remarks.**

- (1) The matrix  $A$  is destroyed.
- (2) CMSLV and CMSLV1 produce the same results for  $X$  and the inverse of  $A$ .

**Algorithm.** The partial pivot Gauss elimination procedure is used. The pivots  $a_{kj}$  are selected so that  $|\operatorname{Re}(a_{kj})| + |\operatorname{Im}(a_{kj})| = \max \{ |\operatorname{Re}(a_{ij})| + |\operatorname{Im}(a_{ij})| : i = j, \dots, n \}$ .

**Programming.** CMSLV1 calls the LINPACK subroutines CGEFA, CGESL, and CGEDI. These subroutines were written by Cleve Moler (University of New Mexico). The subroutines CSWAP, CSCAL, CAXPY and functions CDOTC, ICAMAX are also used.

**Reference.** Dongarra, J. J., Bunch, J. R., Moler, C. B., and Stewart, G. W., *LINPACK Users' Guide*, SIAM, 1979.

**CALL DCMSLV(MO, $n$ , $m$ ,AR,AI, $ka$ ,BR,BI, $kb$ ,IERR,IPVT,WK)**

AR and AI are double precision matrices of order  $n \geq 1$ . AR and AI are the real and imaginary parts of the complex matrix  $A$  whose inverse is to be computed or for which  $AX = B$  is to be solved. If  $MO = 0$  then the inverse of the complex matrix is computed and the results stored in AR and AI. If  $MO \neq 0$  then the inverse is not computed.

The argument  $m$  is an integer. If  $m \geq 1$  then BR and BI are double precision matrices having  $n$  rows and  $m$  columns. In this case, BR and BI are the real and imaginary parts of the complex matrix  $B$  for which  $AX = B$  is to be solved. When DCMSLV is called, the real and imaginary parts of the solution  $X$  are computed and stored in BR and BI. If  $m \leq 0$  then there are no equations to be solved. In this case BR and BI are ignored.

The argument  $ka$  is the number of rows in the dimension statements for AR and AI in the calling program, and the argument  $kb$  is the number of rows in the dimension statements for BR and BI in the calling program. If  $m \leq 0$  then the argument  $kb$  is ignored.

IPVT is an integer array of dimension  $n$  or larger that is used by the routine for keeping track of the row interchanges that are made.

WK is a double precision array of dimension  $2n$  or larger that is a work space for the routine. If  $MO \neq 0$  then WK is ignored.

**Error Return.** IERR is a variable that reports the status of the results. When DCMSLV terminates IERR has one of the following values:

- IERR = 0 The requested results were obtained.
- IERR = -1 Either  $n$ ,  $ka$ , or  $kb$  is incorrect.
- IERR =  $k$  The  $k^{\text{th}}$  columns of AR and AI have been reduced to columns containing zeros. The requested results were not obtained.

**Remark.** The matrices AR and AI are destroyed.

**Algorithm.** The partial pivot Gauss elimination procedure is used. The pivots  $A_{kj}$  are selected so that  $|\operatorname{Re}(a_{kj})| + |\operatorname{Im}(a_{kj})| = \max \{|\operatorname{Re}(a_{ij})| + |\operatorname{Im}(a_{ij})| : i = j, \dots, n\}$ .

**Programming.** DCMSLV calls the subroutines DCFAC, DCSOL, and DCMINV. These routines were written by A. H. Morris. The functions CDIVID and DPMPAR are also used.

## SOLUTION OF COMPLEX EQUATIONS WITH ITERATIVE IMPROVEMENT

Given a complex  $n \times n$  matrix  $A$  and a complex column vector  $b$ . The following routine is available for solving the equation  $Ax = b$ . Iterative improvement is performed to compute the solution  $x$  to machine accuracy.

**CALL CSLVMP(MO,n,A,ka,b,X,WK,IWK,IND)**

MO is an input argument which specifies if CSLVMP is being called for the first time. On an initial call, MO = 0 and we have the following setup:

$A$  is a 2-dimensional complex array of dimension  $ka \times n$  containing the matrix  $A$ ,  $b$  a complex vector of dimension  $n$ , and  $X$  a complex array of dimension  $n$ . When CSLVMP is called,  $Ax = b$  is solved and the solution stored in  $X$ .  $A$  and  $b$  are not modified by the routine.

WK is a complex array of dimension  $n^2 + n$  or larger, and IWK an integer array of dimension  $n$  or larger. These arrays are for internal use by the routine. On an initial call to CSLVMP, an LU decomposition is obtained for  $A$  and stored in WK and IWK. Then the equation  $Ax = b$  is solved.

IND is an integer variable that reports the status of the results. On an initial call to CSLVMP, when the routine terminates IND has one of the following values:

IND = 0 The solution  $X$  was obtained to machine accuracy.

IND = 1  $X$  was obtained, but not to machine accuracy.

IND =  $-k$  The  $k^{th}$  column of  $A$  was reduced to a column containing only zeros. In this case no solution can be obtained.

After an initial call to CSLVMP, if IND = 0 or 1 on output, then the routine may be called to solve a new set of equations  $Ax = \tilde{b}$  without having to redecompose the matrix  $A$ . In this case, the input argument MO may be set to any nonzero value. When MO  $\neq$  0 it is assumed that only  $b$  has been modified. The routine employs the LU decomposition obtained on the initial call to CSLVMP to solve the new set of equations  $Ax = b$ . On output  $X$  will contain the solution to the new set of equations. As before,  $A$  and  $b$  are not modified by the routine.

If CSLVMP is recalled with MO  $\neq$  0, then when the routine terminates IND has one of the following values:

IND = 0 The solution  $X$  was obtained to machine accuracy.

IND = 1  $X$  was obtained, but not to machine accuracy.

**Programming.** CSLVMP calls the subroutine CLUIIMP. These routines were written by A. H. Morris. The subroutines CMCOPY, CGEFA, CGESL, CSCAL, CAXPY and functions SPMPAR, CDOTC, ICAMAX are also employed.

## SINGULAR VALUE DECOMPOSITION OF A MATRIX

If  $A$  is a complex  $m \times n$  matrix then there exists an  $m \times m$  unitary matrix  $U$  and an  $n \times n$  unitary matrix  $V$  such that  $D = U^*AV$  is a diagonal matrix<sup>1</sup>. Let  $d_1, \dots, d_k$  be the diagonal elements of  $D$  where  $k = \min\{m, n\}$ . Then  $U$  and  $V$  can be selected so that the diagonal elements are real numbers and  $d_1 \geq d_2 \geq \dots \geq d_k \geq 0$ . The nonnegative diagonal elements  $d_i$  are unique, and if  $A$  is a real matrix then  $U$  and  $V$  can be chosen to be real orthogonal matrices. The decomposition  $D = U^*AV$  is called the *singular value decomposition of  $A$* . The elements  $d_1, \dots, d_k$  are the *singular values* of  $A$ , the columns of  $U$  are *left singular vectors*, and the columns of  $V$  are *right singular vectors*.

**Remark.** For  $m > n$ ,  $D = \begin{pmatrix} D_1 \\ 0 \end{pmatrix}$  where  $D_1 = \text{diag}(d_1, \dots, d_n)$ . Consequently, if  $U$  is partitioned into  $U = (U_1, U_2)$  where  $U_1$  has  $n$  columns, then it follows that  $A = UDV^* = U_1D_1V^*$ . The decomposition  $A = U_1D_1V^*$  is frequently also called the singular value decomposition, and in many applications it suffices.

The following subroutines are available for finding the singular value decomposition  $D = U^*AV$  of a matrix  $A$ .

```
CALL SSVDC(A, ka, m, n, D, E, U, ku, V, kv, WORK, JOB, INFO)
CALL DSVDC(A, ka, m, n, D, E, U, ku, V, kv, WORK, JOB, INFO)
CALL CSVDC(A, ka, m, n, D, E, U, ku, V, kv, WORK, JOB, INFO)
```

$A$  is a 2-dimensional array of dimension  $ka \times n$  containing the  $m \times n$  matrix whose singular value decomposition is to be computed.  $D$  is an array of dimension  $\min\{m+1, n\}$ . When any of the routines is called, the singular values of  $A$  are computed and stored in descending order of magnitude in  $D(1), \dots, D(k)$  where  $k = \min\{m, n\}$ .

$JOB$  is an integer that controls the computation of the singular vectors. It is assumed that  $JOB = I \cdot 10 + J$  when  $I, J = 0, 1, \dots, 9$ .  $I$  and  $J$  have the following meaning.

- $I = 0$  Do not compute the left singular vectors.
- $I = 1$  Compute all  $m$  left singular vectors.
- $I > 1$  Compute the first  $\min\{m, n\}$  left singular vectors. (Here we compute the decomposition  $A = U_1D_1V^*$ .)
- $J = 0$  Do not compute the right singular vectors.
- $J > 0$  Compute the right singular vectors.

$U$  is a 2-dimensional array which contains the left singular vectors that are requested, and  $ku$  is the number of rows in the dimension statement for  $U$  in the calling program. It is assumed that  $ku \geq m$ . If no left singular vectors are requested (i.e., if  $JOB < 10$ ) then  $U$  is ignored by the routines. Otherwise,  $U$  must be of dimension  $ku \times m$  if all  $m$  left singular vectors are requested, and  $U$  must be of dimension  $ku \times \min\{m, n\}$  if the first  $\min\{m, n\}$  left singular vectors are requested.

<sup>1</sup> $U^*$  denotes the adjoint matrix of  $U$ .

$V$  is a 2-dimensional array which contains the right singular vectors that are requested, and  $kv$  is the number of rows in the dimension statement for  $V$  in the calling program. It is assumed that  $kv \geq n$ . If no right singular vectors are requested then  $V$  is ignored by the routines. Otherwise,  $V$  must be of dimension  $kv \times n$  if the right singular vectors are requested.

$E$  is an array of dimension  $n$  or larger, and  $WORK$  is an array of dimension  $m$  or larger.  $E$  and  $WORK$  are storage areas for the routines.

#### Remarks.

- (1) If SSVDC is called then it is assumed that the arrays  $A, D, E, U, V, WORK$  are real arrays, if DSVDC is called then it is assumed that the arrays are double precision arrays, and if CSVDC is called then it is assumed that the arrays are complex arrays.
- (2) The contents of  $A$  are destroyed by the routines. If left singular vectors are requested and there is sufficient storage in  $A$  to hold the vectors (there will be sufficient storage if  $m \leq n$  or  $JOB \geq 20$ ), then one may set  $U = A$ . Similarly, if right singular vectors are requested and  $m \geq n$  then one may set  $V = A$ . However, only one of the two arrays  $U$  and  $V$  may be identified with  $A$ .

**Error Return.** INFO is an integer variable. If all the singular values are found then INFO will be set to 0 and the array  $E$  will contain zeros. However, if the  $j^{th}$  singular value cannot be found then INFO is set to  $j$ . In this case, if  $j < k$  where  $k = \min\{m, n\}$  then the singular values  $d_{j+1}, \dots, d_k$  will have been computed and stored in  $D$ .  $A$  will have been reduced to an upper bidiagonal matrix  $B$  with  $D$  as its diagonal and  $E$  its super diagonal. If  $U$  and  $V$  have been requested then  $B = U^*AV$  will be satisfied.

**Programming.** SSVDC, DSVDC, and CSVDC are part of the LINPACK package of matrix subroutines released by Argonne National Laboratory. The routines were coded by G. W. Stewart (University of Maryland). The routines employ the vector subroutines SSWAP, SROT, SDOT, SSCAL, SAXPY, SNRM2, DSWAP, DROT, DDOT, DSCAL, DAXPY, DNRM2, and CSWAP, CSROT, CDOTC, CSCAL, CAXPY, SCNRM2. Also the subroutines SROTG and DROTG are called.

**Reference.** Dongarra, J. J., Bunch, J. R., Moler, C. B., and Stewart, G. W., *LINPACK Users' Guide*, Society for Industrial and Applied Mathematics, Philadelphia, 1979.

## EVALUATION OF THE CHARACTERISTIC POLYNOMIAL OF A MATRIX

The following functions are available for computing the determinant of  $A - xI$  where  $A$  is an  $n \times n$  matrix,  $x$  a number, and  $I$  the  $n \times n$  identity matrix.

**DET**( $A, ka, n, x$ )  
**DPDET**( $A, ka, n, x$ )  
**CDET**( $A, ka, n, x$ )

**DET** is a real function that is used when  $A$  is a real matrix and  $x$  a real number, **DPDET** is a double precision function that is used when  $A$  is a double precision matrix and  $x$  a double precision number, and **CDET** is a complex function that is used when  $A$  is a complex matrix and  $x$  a complex number.

The value of the appropriate function is the determinant of the matrix  $A - xI$ . The argument  $ka$  has the value:

$ka$  = the number of rows in the dimension statement for  $A$  in the calling program  
It is assumed that  $ka \geq n \geq 1$ .

**Note.**  $A$  is destroyed during computation.

**Algorithm.** Gauss partial pivoting is performed to reduce  $A - xI$  to upper triangular form. In **CDET** the pivots  $a_{kj}$  are selected so that  $|\operatorname{Re}(a_{kj})| + |\operatorname{Im}(a_{kj})| = \operatorname{Max}\{|\operatorname{Re}(a_{ij})| + |\operatorname{Im}(a_{ij})| : i = j, \dots, n\}$  rather than  $|a_{kj}| = \max\{|a_{ij}| : i = j, \dots, n\}$ .

**Programmer.** A. H. Morris

## SOLUTION OF THE MATRIX EQUATION $AX + XB = C$

Given an  $m \times m$  matrix  $A$ ,  $n \times n$  matrix  $B$ , and  $m \times n$  matrix  $C$ . The subroutines ABSLV and DABSLV are available for obtaining the  $m \times n$  matrix  $X$  which solves the equation  $AX + XB = C$ . ABSLV yields single precision results and DABSLV yields double precision results.

**CALL ABSLV(MO,m,n,A,ka,B,kb,C,kc,WK,IND)**  
**CALL DABSLV(MO,m,n,A,ka,B,kb,C,kc,WK,IND)**

If ABSLV is called then it is assumed that  $A, B, C$ , and  $WK$  are real arrays. Otherwise, if DABSLV is called then it is assumed that  $A, B, C$ , and  $WK$  are double precision arrays.

It is assumed that  $m \geq 1$  and  $n \geq 1$ . The input arguments  $ka, kb, kc$  have the following values:

$ka$  = the number of rows in the dimension statement for  $A$  in the calling program

$kb$  = the number of rows in the dimension statement for  $B$  in the calling program

$kc$  = the number of rows in the dimension statement for  $C$  in the calling program

It is required that  $ka \geq m, kb \geq n, kc \geq m$ .

$WK$  is an array of dimension  $m^2 + n^2 + 2k$  or larger where  $k = \max\{m, n\}$ .  $WK$  is a general storage area for the routine.

$MO$  is an input argument which specifies if the routine is being called for the first time. On an initial call  $MO = 0$ . In this case,  $A$  is reduced to lower real Schur form,  $B$  is reduced to upper real Schur form, and then the transformed system of equations is solved.

$IND$  is a variable that reports the status of the results. When the routine terminates,  $IND$  has one of the following values:

$IND = 0$  The solution was obtained and stored in  $C$ .

$IND = 1$  The equations are inconsistent for  $A$  and  $B$ .

$IND = -1$   $A$  could not be reduced to lower Schur form.

$IND = -2$   $B$  could not be reduced to upper Schur form.

If  $IND \neq 0$  then no solution is obtained.

When  $IND = 0$ ,  $A$  contains the lower Schur form of the matrix  $A$ ,  $B$  contains the upper Schur form of the matrix  $B$ , and  $WK$  contains the orthogonal matrices involved in the decompositions of  $A$  and  $B$ . This information can be reused to solve a new set of equations. The following options are available:

$MO = 1$  New matrices  $A$  and  $C$  are given. The data for  $B$  is reused in solving the new set of equations.

$MO = 2$  New matrices  $B$  and  $C$  are given. The data for  $A$  is reused in solving the new set of equations.

$MO \neq 0, 1, 2$  A new matrix  $C$  is given. The data for  $A$  and  $B$  is reused in solving the new set of equations.

When the routine is recalled, it is assumed that  $m$ ,  $n$ , and  $WK$  have not been modified.

**Programming.** ABSLV employs the subroutines ABSLV1, ORTHES, ORTRN1, SCHUR, SHRSLV, SLV, and DABSLV employs the routines DABSV1, DORTH, DRTRN1, DSCHUR, DSHSLV, DPSLV. ABSLV and DABSLV are adaptations by A. H. Morris of the subroutine AXPXB written by R. H. Bartels and G. W. Stewart (University of Texas at Austin).

**Reference.** Bartels, R. H. and Stewart, G. W., "Algorithm 432, Solution of the Matrix Equation  $AX + XB = C$ ," *Comm. ACM* 15 (1972), pp. 820-826.



## SOLUTION OF THE MATRIX EQUATION $A^tX + XA = C$ WHEN $C$ IS SYMMETRIC

Given matrices  $A$  and  $C$  of order  $n$  where  $C$  is symmetric. Then the subroutines TASLV and DTASLV are available for obtaining the symmetric matrix  $X$  which solves the equation  $A^tX + XA = C$ . TASLV yields single precision results and DTASLV yields double precision results.

CALL TASLV(MO,n, A, ka, C, kc, WK, IND)  
CALL DTASLV(MO,n, A, ka, C, kc, WK, IND)

If TASLV is called then it is assumed that  $A$ ,  $C$  and  $WK$  are real arrays. Otherwise, if DTASLV is called then it is assumed that  $A$ ,  $C$ , and  $WK$  are double precision arrays.

It is assumed that  $n \geq 1$ . The input arguments  $ka$  and  $kc$  have the following values:

$ka$  = the number of rows in the dimension statement for  $A$  in the calling program

$kc$  = the number of rows in the dimension statement for  $C$  in the calling program

It is required that  $ka \geq n$  and  $kc \geq n$ .

$WK$  is an array of dimension  $n^2 + 2n$  or larger that is a general storage area for the routine.

$MO$  is an input argument which specifies if the routine is being called for the first time. On an initial call  $MO = 0$ . In this case,  $A$  is reduced to upper real Schur form and then the transformed system of equations is solved.

$IND$  is a variable that reports the status of the results. When the routine terminates,  $IND$  has one of the following values:

$IND = 0$  The solution was obtained and stored in  $C$ .

$IND = 1$  The equations are inconsistent.

$IND = -1$   $A$  could not be reduced to upper Schur form.

If  $IND \neq 0$  then no solution is obtained.

When  $IND = 0$ ,  $A$  contains the upper Schur form of the matrix  $A$  and  $WK$  contains the orthogonal matrix involved in the decomposition of  $A$ . This data can be reused to solve a new set of equations  $A^tX + XA = \tilde{C}$ . In this case,  $MO$  can be set to any nonzero value. When  $MO \neq 0$  it is assumed that only  $C$  has been modified. When the routine terminates, the solution for the new set of equations is stored in  $C$ .

**Programming.** TASLV employs the subroutines TASLV1, ORTHES, ORTRN1, SCHUR, SYMSLV, SLV and DTASLV employs the routines DTASV1, DORTH, DRTRN1, DSCHUR, DSYSV, DPSLV. TASLV and DTASLV are adaptations by A. H. Morris of the subroutine ATXPXA written by R. H. Bartels and G. W. Stewart (University of Texas at Austin).

**Reference.** Bartels, R. H. and Stewart, G. W., "Algorithm 432, Solution of the Matrix Equation  $AX + XB = C$ ," *Comm. ACM* 15 (1972), pp. 820-826.

## SOLUTION OF THE MATRIX EQUATION $AX^2 + BX + C = 0$

Given complex  $n \times n$  matrices  $A, B$ , and  $C$ . The following subroutine is available for obtaining a complex  $n \times n$  matrix  $X$  which solves the equation  $AX^2 + BX + C = 0$ .

**CALL SQUINT**( $m, n, A, B, C, \text{IND}, X, \text{WK}, \ell, \tau, \text{MAX}, \text{IERR}$ )

It is assumed that  $A, B, C$ , and  $X$  are 2-dimensional complex arrays of dimension  $m \times n$  where  $m \geq n$ . When SQUINT is called, the  $n \times n$  complex matrix solution obtained for  $AX^2 + BX + C = 0$  is stored in  $X$ .  $A, B$ , and  $C$  are modified by the routine.

IND is an integer variable. On input, if  $\text{IND} \neq 0$  then it is assumed that an initial approximation for the desired solution is provided in  $X$  by the user. Otherwise, if  $\text{IND} = 0$  then the routine provides its own initial approximation. Then Newton iteration is performed. On output,  $\text{IND} =$  the number of iterations that were performed to compute  $X$ .

WK is a complex array of dimension  $\ell$  that is a work space for the routine. It is assumed that  $\ell \geq 7n^2 + n$ . When SQUINT terminates,  $\text{WK}(1)$  is a complex number whose real part is the norm  $\|AX^2 + BX + C\|_\infty$ .

The argument  $\tau$  is a real number. If  $\tau \leq 0$  then  $X$  is computed to machine precision. Otherwise, if  $\tau > 0$  then iteration terminates when  $\|AX^2 + BX + C\| < \tau$ .

MAX is a variable. If  $\text{MAX} > 0$  then MAX is the maximum number of iterations that may be performed. If  $\text{MAX} \leq 0$  then it is reset by the routine to 30, the default maximum number of iterations.

**Error Return.** IERR is a variable that is set by the routine. If a solution  $X$  is obtained, then IERR is assigned the value 0. Otherwise, IERR has one of the following values:

- IERR = 1      MAX iterations were performed. More iterations are needed.
- IERR = 2, 3    Factorization of the equations could not be completed.  $X$  cannot be computed.
- IERR = 10 +  $n$  Newton iteration failed on iteration  $n$ . Possibly too much accuracy was requested.  $X$  cannot be computed.
- IERR = 999    (Input Error) Either  $n < 1$ ,  $m < n$ , or  $\ell < 7n^2 + n$ .

When IERR = 1 occurs,  $X$  contains the most recent value obtained for the solution and  $\text{WK}(1)$  is a complex number whose real part is the latest value obtained for the norm  $\|AX^2 + BX + C\|_\infty$ .

**Programming.** SQUINT employs the subroutines SQUIN2, CQZHES, CQZIT, TRISLV, and CTRANS. These routines were designed by George J. Davis (Georgia State University, Atlanta, Georgia). CQZHES and CQZIT are modifications of the EISPACK subroutines QZHES and QZIT, developed at Argonne National Laboratory. The function SPMPAR is also used.

#### References.

- (1) Davis, G. J., "Algorithm 598. An Algorithm to Compute Solvents of the Matrix Equation  $AX^2 + BX + C = 0$ ," *ACM Trans. Math Software* 9 (1983), pp. 246-254.
- (2) Garbow, B. S., et al., *Matriz Eigensystems Routines - EISPACK Guide Extension*, Springer-Verlag, 1977.

## EXPONENTIAL OF A REAL MATRIX

Let  $A$  be a real matrix of order  $n \geq 1$ . Then the subroutines MEXP and DMEXP are available for computing the exponential matrix  $e^A = \sum_{i=0}^{\infty} A^i/i!$ . MEXP yields single precision results and DMEXP yields double precision results.

**CALL MEXP**( $A, ka, n, Z, kz, WK, IERR$ )

$A$  is a real matrix of order  $n \geq 1$  and  $Z$  a real 2-dimensional array. MEXP computes  $e^A$  and stores the results in  $Z$ . The arguments  $ka$  and  $kz$  have the following values:

$ka$  = the number of rows in the dimension statement for  $A$  in the calling program

$kz$  = the number of rows in the dimension statement for  $Z$  in the calling program

It is assumed that  $ka \geq n$ ,  $kz \geq n$ , and that  $A$  and  $Z$  are different storage areas. If  $n > 1$  then  $A$  is destroyed.

$WK$  is a real array of dimension  $n(n+8)$  or larger that is a work space for the routine.

$IERR$  is a variable that reports the status of the results. When MEXP terminates,  $IERR$  is assigned one of the following values:

$IERR = 0$  The exponential was computed.

$IERR = 1$  The norm  $\|A\|_1 = \max_j \sum_i |a_{ij}|$  is too large.  $e^A$  cannot be computed.

$IERR = 2$  The Pade denominator matrix is singular. (This should never occur.)

**Algorithm.**  $A$  is balanced, yielding a matrix  $B = D^{-1}P^tAPD$  where  $D$  is a diagonal matrix,  $P$  a permutation matrix, and  $\|B\|_1 \leq \|A\|_1$ . Next  $m$  is set to the smallest nonnegative integer such that  $\min\{\|B\|_1, \|B\|_\infty\} \leq 2^m$ , and the 8<sup>th</sup> diagonal Pade approximation for  $e^x$  is used to compute  $\exp(B/2^m)$ . Then  $e^B = [\exp(B/2^m)]^{2^m}$  is obtained by  $m$  squarings, and  $e^A = PDe^B D^{-1}P^t$  is applied.

**Programming.** MEXP calls the subroutines BALANC, BALINV, and SLV. The function IPMPAR is also used. MEXP was written by A. H. Morris.

**Reference.** Ward, Robert, C., "Numerical Computation of the Matrix Exponential with Accuracy Estimate," *SIAM J. Numerical Analysis* 14 (1977), pp. 600-610

**CALL DMEXP**( $A, ka, n, Z, kz, WK, IERR$ )

$A$  is a double precision matrix of order  $n \geq 1$  and  $Z$  a double precision 2-dimensional array. DMEXP computes  $e^A$  and stores the results in  $Z$ . The arguments  $ka$  and  $kz$  have the following values:

$ka$  = the number of rows in the dimension statement for  $A$  in the calling program

$kz$  = the number of rows in the dimension statement for  $Z$  in the calling program

It is assumed that  $ka \geq n$ ,  $kz \geq n$ , and that  $A$  and  $Z$  are different storage areas. If  $n > 1$  then  $A$  is destroyed.

WK is a double precision array of dimension  $n(n + 12)$  or larger that is a work space for the routine.

IERR is a variable that reports the status of the results. When DMEXP terminates, IERR is assigned one of the following values:

IERR = 0 The exponential was computed.

IERR = 1 The norm  $\|A\|_1$  is too large.  $e^A$  cannot be computed.

IERR = 2 The Pade denominator matrix is singular. (This should never occur.)

**Programming.** DMEXP calls the subroutines DBAL, DBALNV, and DPSLV. The function IPMPAR is also used. DMEXP was written by A. H. Morris.

**Reference.** Ward, Robert C., "Numerical Computation of the Matrix Exponential with Accuracy Estimate," *SIAM J. Numerical Analysis* 14 (1977), pp. 600-610

## SOLVING SYSTEMS OF 200-400 LINEAR EQUATIONS

For  $n \geq 1$ , let  $A$  denote an  $n \times n$  matrix and  $b$  a column vector of dimension  $n$ . Then the subroutines LE, DPLE, and CLE are available for solving the equations  $Ax = b$  where  $A$  is not stored in-core. For large  $n$ , these routines require a work space of dimension  $\approx n^2/4$ . This permits the solution of systems of equations of double the order permitted by the standard solution procedures.

```
CALL LE(ROWK,n,b,X,WK,IWK,IERR)
CALL DPLE(ROWK,n,b,X,WK,IWK,IERR)
CALL CLE(ROWK,n,b,X,WK,IWK,IERR)
```

$X$  is an array of dimension  $n$  and  $IERR$  an integer variable. When the equations are solved, then  $IERR$  is set to 0 and the solution is stored in  $X$ .

ROWK is the name of a user defined subroutine that has the format:

```
CALL ROWK(n,k,R)
```

$R$  is an array of dimension  $n$  and  $k = 1, \dots, n$ . When ROWK is called, the  $k^{th}$  row of the matrix  $A$  is stored in  $R$ . ROWK must be declared in the calling program to be of type EXTERNAL.

WK is an array of dimension  $[n^2/4] + n + 3$  or larger,<sup>1</sup> and IWK is an integer array of dimension  $\max\{1, n - 1\}$  or larger. WK and IWK are work spaces for the routines.

**Error Return.**  $IERR = k$  when the first  $k$  rows of  $A$  are found to be linearly dependent.

### Remarks.

- (1) When LE is called then it is assumed that  $b$ ,  $X$ , WK and the array  $R$  in ROWK are real arrays. When DPLE is called then it is assumed that these arrays are double precision arrays, and when CLE is called then it is assumed that the arrays are complex.
- (2) When the equations are solved, ROWK is called to attach the first row of  $A$ , then the second row, etc. Each row of  $A$  is attached only once.
- (3) The array  $b$  is not modified by the routines.

**Example.** Consider a system of  $n = 300$  real linear equations  $Ax = b$  where the rows of  $A$  are stored, one row per logical record, in sequence in an unformatted file (say file 4). Then the following code can be used to solve the equations:

```
REAL B(300),X(300),WK(22803)
INTEGER IWK(300)
EXTERNAL GETROW
DATA N/300/
```

<sup>1</sup>Here  $[n^2/4]$  denotes the largest integer  $\leq n^2/4$ .

```
.  
. .  
REWIND 4  
CALL LE(GETROW,N,B,X,WK,IWK,IERR)
```

Here GETROW may be defined by:

```
SUBROUTINE GETROW(N,I,R)  
REAL R(N)  
READ(4) (R(J),J=1,N)  
RETURN  
END
```

**Algorithm.** The partial pivot Henderson-Wassying procedure is used.

**Programming.** LE, DPLE, and CLE are modified versions (by A. H. Morris) of the subroutine TE, written by A. Wassying (University of the Witwatersrand, Johannesburg, South Africa).

**Reference.** Wassying, A., "Solving  $Ax = b$ : A Method with Reduced Storage Requirements," *SIAM J. Numerical Analysis* 19 (1982), pp. 197-204.

## BAND MATRIX STORAGE

For an  $m \times n$  matrix  $A = (a_{ij})$ , let  $m_\ell$  be the number of diagonals below the main diagonal containing nonzero elements, and  $m_u$  the number of diagonals above the main diagonal containing nonzero elements. Then  $m_\ell$  and  $m_u$  are called the *lower* and *upper band widths* of  $A$ , and  $m_\ell + m_u + 1$  the *total band width* of  $A$ . It is clear that  $0 \leq m_\ell < m$  and  $0 \leq m_u < n$ , and that  $a_{ij} \neq 0$  only when  $i - m_\ell \leq j \leq i + m_u$ . If the band width  $m_\ell + m_u + 1$  is sufficiently small, then it is also clear that a considerable savings in storage can occur by storing only the nonzero diagonals of  $A$ . The band storage scheme adopted by the NSWCL library is to store  $A$  as an  $m \times (m_\ell + m_u + 1)$  matrix  $B = (b_{ik})$ . The columns of  $B$  are the nonzero diagonals of  $A$ . Specifically, for each nonzero  $a_{ij}$ ,  $b_{ik} = a_{ij}$  where  $k = j - i + m_\ell + 1$ . The remaining  $b_{ik}$ 's are zeros.

**Example.** Consider the matrix

$$A = \begin{pmatrix} a_{11} & a_{12} & a_{13} & 0 & 0 & 0 & 0 \\ a_{21} & a_{22} & a_{23} & a_{24} & 0 & 0 & 0 \\ 0 & a_{32} & a_{33} & a_{34} & a_{35} & 0 & 0 \\ 0 & 0 & a_{43} & a_{44} & a_{45} & a_{46} & 0 \\ 0 & 0 & 0 & a_{54} & a_{55} & a_{56} & a_{57} \\ 0 & 0 & 0 & 0 & a_{65} & a_{66} & a_{67} \\ 0 & 0 & 0 & 0 & 0 & a_{76} & a_{77} \\ 0 & 0 & 0 & 0 & 0 & 0 & a_{87} \end{pmatrix}$$

where  $m_\ell = 1$  and  $m_u = 2$ . Then  $A$  will be stored in band form as follows:

$$B = \begin{pmatrix} 0 & a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} & a_{24} \\ a_{32} & a_{33} & a_{34} & a_{35} \\ a_{43} & a_{44} & a_{45} & a_{46} \\ a_{54} & a_{55} & a_{56} & a_{57} \\ a_{65} & a_{66} & a_{67} & 0 \\ a_{76} & a_{77} & 0 & 0 \\ a_{87} & 0 & 0 & 0 \end{pmatrix}$$

**Remark.** The first  $m_\ell$  columns of  $B$  contain the nonzero diagonals of  $A$  below the main diagonal, the  $(m_\ell + 1)^{st}$  column of  $B$  contains the main diagonal, and the last  $m_u$  columns of  $B$  contain the nonzero diagonals of  $A$  above the main diagonal.



## CONVERSION OF BANDED MATRICES TO AND FROM THE STANDARD FORMAT

The following subroutines permit one to convert matrices to and from the standard format.

**CALL CVBR( $A, ka, m, n, m_\ell, m_u, B, kb$ )**  
**CALL CVBC( $A, ka, m, n, m_\ell, m_u, B, kb$ )**

$A$  is an  $m \times n$  matrix stored in band form,  $m_\ell$  the number of diagonals below the main diagonal containing nonzero elements, and  $m_u$  the number of diagonals above the main diagonal containing nonzero elements. The argument  $ka$  is the number of rows in the dimension statement for  $A$  in the calling program. It is assumed that  $0 \leq m_\ell < m, 0 \leq m_u < n$ , and  $ka \geq m$ .

$B$  is a 2-dimensional array of dimension  $kb \times n$  where  $kb \geq m$ . CVBR is used if  $A$  and  $B$  are real arrays, and CVBC is used if  $A$  and  $B$  are complex arrays. When the routine is called, the matrix  $A$  is stored in the array  $B$  in the standard format.

**Remark.**  $B$  may begin in the same location as  $A$ . If  $B$  begins in the same location then it is assumed that  $kb = ka$ . In this case, the result  $B$  will overwrite the input data  $A$ . Otherwise, if  $B$  does not begin in the same location as  $A$ , then it is assumed that the storage areas  $A$  and  $B$  do not overlap.

**Programmer.** A. H. Morris

**CALL CVRB( $A, ka, m, n, m_\ell, m_u, B, kb$ )**  
**CALL CVCB( $A, ka, m, n, m_\ell, m_u, B, kb$ )**

$A$  is an  $m \times n$  matrix stored in the standard format, and  $m_\ell$  and  $m_u$  are integers such that  $0 \leq m_\ell < m$  and  $0 \leq m_u < n$ . The argument  $ka$  is the number of rows in the dimension statement for  $A$  in the calling program. It is assumed that  $ka \geq m$ .

$B$  is a 2-dimension array of dimension  $kb \times \ell$  where  $kb \geq m$  and  $\ell \geq m_\ell + m_u + 1$ . CVRB is used if  $A$  and  $B$  are real arrays, and CVCB is used if  $A$  and  $B$  are complex arrays. When the routine is called, the  $m_\ell$  diagonals of  $A$  immediately below the main diagonal, the main diagonal, and the  $m_u$  diagonals immediately above the main diagonal are stored in band form in  $B$ .

### Remarks.

- (1) Given a matrix  $A = (a_{ij})$ , then these routines may be used to convert  $A$  to band form when the lower and upper bandwidths  $m_\ell$  and  $m_u$  of  $A$  are known. If  $m_\ell$  and  $m_u$  are not known, then the subroutines CVRB1 and CVCB1 described below can be used to convert  $A$  to band form.
- (2)  $B$  may begin in the same location as  $A$ . If  $B$  begins in the same location then it is assumed that  $kb = ka$ . In this case, the result  $B$  will overwrite the input data  $A$ .

Otherwise, if  $B$  does not begin in the same location as  $A$ , then it is assumed that the storage areas  $A$  and  $B$  do not overlap.

**Programmer.** A. H. Morris

```
CALL CVRB1(A, ka, m, n, ml, mu, B, kb,  $\ell$ , IERR)  
CALL CVCB1(A, ka, m, n, ml, mu, B, kb,  $\ell$ , IERR)
```

$A$  is an  $m \times n$  matrix stored in the standard format. The argument  $ka$  is the number of rows in the dimension statement for  $A$  in the calling program. It is assumed that  $A$  is to be stored in band form in  $B$ .  $B$  is a 2-dimensional array of dimension  $kb \times \ell$  where  $kb \geq m$ . The argument  $\ell$  is an estimate of the maximum number of diagonals of  $A$  that will have to be stored.

CVRB1 is used if  $A$  and  $B$  are real arrays, and CVCB1 is used if  $A$  and  $B$  are complex arrays. IERR,  $m_l$ , and  $m_u$  are integer variables. When the routine is called, if  $\ell$  specifies sufficient storage for  $B$  then  $A$  is stored in band form in  $B$ . Also IERR is assigned the value 0,  $m_l$  = the number of diagonals of  $A$  below the main diagonal containing nonzero elements, and  $m_u$  = the number of diagonals above the main diagonal containing nonzero elements.

**Error Return.** If  $\ell$  does not specify sufficient storage for  $B$ , then IERR is assigned the value  $m_l + m_u + 1$ . Reset  $\ell \geq$  IERR.

**Remark.**  $B$  may begin in the same location as  $A$ . If  $B$  begins in the same location then it is assumed that  $kb = ka$ . In this case, the result  $B$  will overwrite the input data  $A$ . Otherwise, if  $B$  does not begin in the same location as  $A$ , then it is assumed that the storage areas  $A$  and  $B$  do not overlap.

**Programming.** CVRB1 calls the subroutine CVRB, and CVCB1 calls the subroutine CVCB. These routines were written by A. H. Morris.

## CONVERSION OF BANDED MATRICES TO AND FROM SPARSE FORM

The following subroutines permit one to convert matrices to and from sparse form.

```
CALL MCVBS(A,ka,m,n,ml,mu,B,IB,JB,NUM,IERR)
CALL CMCVBS(A,ka,m,n,ml,mu,B,IB,JB,NUM,IERR)
```

$A$  is an  $m \times n$  matrix stored in band form,  $m_l$  the number of diagonals below the main diagonal containing nonzero elements, and  $m_u$  the number of diagonals above the main diagonal containing nonzero elements. The argument  $ka$  is the number of rows in the dimension statement for  $A$  in the calling program. It is assumed that  $0 \leq m_l < m, 0 \leq m_u < n$ , and  $ka \geq m$ .

It is assumed that  $A$  is to be stored in sparse form in the arrays  $B, IB, JB$ .  $NUM$  is the estimated maximum number of elements that will appear in  $B$  and  $JB$ . It is assumed that  $B$  and  $JB$  are of dimension  $\max\{1, NUM\}$  and that  $IB$  is of dimension  $m + 1$ .

MCVBS is used if  $A$  and  $B$  are real arrays, and CMCVBS is used if  $A$  and  $B$  are complex arrays.  $IERR$  is an integer variable. When the routine is called, if  $NUM$  specifies sufficient storage for  $B$  and  $JB$ , then  $IERR$  is assigned the value 0 and  $A$  is stored in sparse form in  $B, IB, JB$ .

**Error Return.** If there is not sufficient storage in  $B$  and  $JB$  for the  $i^{th}$  row of  $A$ , then  $IERR$  is set to  $i$  and the routine terminates. In this case, if  $i > 1$  then the first  $i - 1$  rows of  $A$  will have been stored in  $B$  and  $JB$ . Also  $IB(1), \dots, IB(i)$  will contain the appropriate row locations.

```
CALL MCVSB(A,IA,JA,m,n,B,kb,l,ml,mu,IERR)
CALL CMCVSB(A,IA,JA,m,n,B,kb,l,ml,mu,IERR)
```

$A$  is an  $m \times n$  sparse matrix stored in the arrays  $A, IA, JA$ . It is assumed that  $A$  is to be stored in band form in  $B$ .  $B$  is a 2-dimensional array of dimension  $kb \times l$  where  $kb \geq m$ . The argument  $l$  is an estimate of the maximum number of diagonals of  $A$  that will have to be stored.

MCVSB is used if  $A$  and  $B$  are real arrays, and CMCVSB is used if  $A$  and  $B$  are complex arrays.  $IERR, m_l$ , and  $m_u$  are integer variables. When the routine is called, if  $l$  specifies sufficient storage for  $B$  then  $A$  is stored in band form in  $B$ . Also  $IERR$  is assigned the value 0,  $m_l$  = the number of diagonals of  $A$  below the main diagonal containing nonzero elements, and  $m_u$  = the number of diagonals above the main diagonal containing nonzero elements.

**Error Return.** If  $l$  does not specify sufficient storage for  $B$ , then  $IERR$  is assigned the value  $m_l + m_u + 1$ . Reset  $l \geq IERR$ .

**Programmer.** A. H. Morris

## TRANSPOSING BANDED MATRICES

The following subroutines are available for transposing banded matrices.

**CALL BPOSE( $A, ka, m, n, m_\ell, m_u, B, kb$ )**  
**CALL CBPOSE( $A, ka, m, n, m_\ell, m_u, B, kb$ )**

$A$  is an  $m \times n$  matrix stored in band form,  $m_\ell$  the number of diagonals below the main diagonal containing nonzero elements, and  $m_u$  the number of diagonals above the main diagonal containing nonzero elements. The argument  $ka$  is the number of rows in the dimension statement for  $A$  in the calling program. It is assumed that  $0 \leq m_\ell < m, 0 \leq m_u < n$ , and  $ka \geq m$ .

$B$  is a 2-dimensional array of dimension  $kb \times \ell$  where  $kb \geq n$  and  $\ell \geq m_\ell + m_u + 1$ . BPOSE is used if  $A$  and  $B$  are real arrays, and CBPOSE is used if  $A$  and  $B$  are complex arrays. When the routine is called, the transpose  $A^t$  of  $A$  is stored in band form in  $B$ .

**Note.** It is assumed that the storage areas  $A$  and  $B$  do not overlap.

**Programmer.** A. H. Morris.

## ADDITION OF BANDED MATRICES

Let  $A$  and  $B$  be  $m \times n$  matrices stored in band form. The following subroutines are available for computing the sum  $C = A + B$ .

**CALL BADD( $m, n, A, ka, m_\ell, m_u, B, kb, n_\ell, n_u, C, kc, \ell, \nu_\ell, \nu_u, IERR$ )**  
**CALL CBADD( $m, n, A, ka, m_\ell, m_u, B, kb, n_\ell, n_u, C, kc, \ell, \nu_\ell, \nu_u, IERR$ )**

$A$  and  $B$  are  $m \times n$  matrices stored in band form,  $m_\ell$  the number of diagonals of  $A$  below the main diagonal containing nonzero elements,  $m_u$  the number of diagonals of  $A$  above the main diagonal containing nonzero elements,  $n_\ell$  the number of diagonals of  $B$  below the main diagonal containing nonzero elements, and  $n_u$  the number of diagonals of  $B$  above the main diagonal containing nonzero elements. The argument  $ka$  is the number of rows in the dimension statement for  $A$  in the calling program, and  $kb$  the number of rows in the dimension statement for  $B$  in the calling program.

It is assumed that  $A + B$  is to be stored in band form in  $C$ .  $C$  is a 2-dimensional array of dimension  $kc \times \ell$  where  $kc \geq m$ . The input argument  $\ell$  is an estimate of the maximum number of diagonals of  $A + B$  which will have to be stored ( $\ell \leq \max\{m_\ell, n_\ell\} + \max\{m_u, n_u\} + 1$ ). BADD is used if  $A$  and  $B$  are real arrays, and CBADD is used if  $A$  and  $B$  are complex arrays. IERR,  $\nu_\ell$ , and  $\nu_u$  are integer variables. When the routine is called, if  $\ell$  specifies sufficient storage for  $C$  then  $A + B$  is computed and stored in band form in  $C$ . Also IERR is assigned the value 0,  $\nu_\ell$  = the number of diagonals of  $A + B$  below the main diagonal containing nonzero elements, and  $\nu_u$  = the number of diagonals of  $A + B$  above the main diagonal containing nonzero elements.

**Error Return.** If  $\ell$  does not specify sufficient storage for  $C$ , then IERR is assigned the value  $\nu$  where  $\nu$  is an estimate of the number of columns needed for  $C$ . Reset  $\ell \geq \nu$ .

**Remarks.** If  $m_\ell \geq n_\ell$  then  $C$  may begin in the same location as  $A$ . If  $C$  begins in the same location as  $A$ , then it is assumed that  $kc = ka$  and that the arrays  $A$  and  $B$  do not overlap. In this case, the result  $C$  will overwrite the input data  $A$ . Similarly, if  $m_\ell \leq n_\ell$  then  $C$  may begin in the same location as  $B$  when  $kc = kb$  and  $A$  and  $B$  do not overlap. Otherwise, if  $C$  does not begin in the same location as  $A$  or  $B$ , then it is assumed that the storage area for  $C$  does not overlap with the storage areas for  $A$  and  $B$ . In this case there is no restriction on  $kc$  (other than the customary restriction that  $kc \geq m$ ).

**Example.** If  $B = -A$  then  $\ell$  may be assigned any value  $\geq 1$ . In this case,  $C$  will contain only the main diagonal of the zero matrix  $A + B$ , and  $\nu_\ell = \nu_u = 0$ .

**Programmer.** A. H. Morris

## SUBTRACTION OF BANDED MATRICES

Let  $A$  and  $B$  be  $m \times n$  matrices stored in band form. The following subroutines are available for computing the difference  $C = A - B$ .

CALL BSUBT( $m, n, A, ka, m_\ell, m_u, B, kb, n_\ell, n_u, C, kc, \ell, \nu_\ell, \nu_u, IERR$ )  
CALL CBSUBT( $m, n, A, ka, m_\ell, m_u, B, kb, n_\ell, n_u, C, kc, \ell, \nu_\ell, \nu_u, IERR$ )

$A$  and  $B$  are  $m \times n$  matrices stored in band form,  $m_\ell$  the number of diagonals of  $A$  below the main diagonal containing nonzero elements,  $m_u$  the number of diagonals of  $A$  above the main diagonal containing nonzero elements,  $n_\ell$  the number of diagonals of  $B$  below the main diagonal containing nonzero elements, and  $n_u$  the number of diagonals of  $B$  above the main diagonal containing nonzero elements. The argument  $ka$  is the number of rows in the dimension statement for  $A$  in the calling program, and  $kb$  the number of rows in the dimension statement for  $B$  in the calling program.

It is assumed that  $A - B$  is to be stored in band form in  $C$ .  $C$  is a 2-dimensional array of dimension  $kc \times \ell$  where  $kc \geq m$ . The input argument  $\ell$  is an estimate of the maximum number of diagonals of  $A - B$  which will have to be stored ( $\ell \leq \max\{m_\ell, n_\ell\} + \max\{m_u, n_u\} + 1$ ). BSUBT is used if  $A$  and  $B$  are real arrays, and CBSUBT is used if  $A$  and  $B$  are complex arrays. IERR,  $\nu_\ell$ , and  $\nu_u$  are integer variables. When the routine is called, if  $\ell$  specifies sufficient storage for  $C$  then  $A - B$  is computed and stored in band form in  $C$ . Also IERR is assigned the value 0,  $\nu_\ell$  = the number of diagonals of  $A - B$  below the main diagonal containing nonzero elements, and  $\nu_u$  = the number of diagonals of  $A - B$  above the main diagonal containing nonzero elements.

**Error Return.** If  $\ell$  does not specify sufficient storage for  $C$ , then IERR is assigned the value  $\nu$  where  $\nu$  is an estimate of the number of columns needed for  $C$ . Reset  $\ell \geq \nu$ .

**Remarks.** If  $m_\ell \geq n_\ell$  then  $C$  may begin in the same location as  $A$ . If  $C$  begins in the same location as  $A$ , then it is assumed that  $kc = ka$  and that the arrays  $A$  and  $B$  do not overlap. In this case, the result  $C$  will overwrite the input data  $A$ . Similarly, if  $m_\ell \leq n_\ell$  then  $C$  may begin in the same location as  $B$  when  $kc = kb$  and  $A$  and  $B$  do not overlap. Otherwise, if  $C$  does not begin in the same location as  $A$  or  $B$ , then it is assumed that the storage area for  $C$  does not overlap with the storage areas for  $A$  and  $B$ . In this case there is no restriction on  $kc$  (other than the customary restriction that  $kc \geq m$ ).

**Example.** If  $B = A$  then  $\ell$  may be assigned any value  $\geq 1$ . In this case,  $C$  will contain only the main diagonal of the zero matrix  $A - B$ , and  $\nu_\ell = \nu_u = 0$ .

**Programmer.** A. H. Morris

## MULTIPLICATION OF BANDED MATRICES

Let  $A$  and  $B$  be matrices stored in band form. The following subroutines are available for computing the product  $C = AB$ .

**CALL BPROD**( $m, n, \ell, A, ka, m_\ell, m_u, B, kb, n_\ell, n_u, C, kc, nc, \nu_\ell, \nu_u, IERR$ )  
**CALL CBPROD**( $m, n, \ell, A, ka, m_\ell, m_u, B, kb, n_\ell, n_u, C, kc, nc, \nu_\ell, \nu_u, IERR$ )

$A$  is an  $m \times n$  matrix stored in band form,  $m_\ell$  the number of diagonals of  $A$  below the main diagonal containing nonzero elements, and  $m_u$  the number of diagonals of  $A$  above the main diagonal containing nonzero elements.  $B$  is an  $n \times \ell$  matrix stored in band form,  $n_\ell$  the number of diagonals of  $B$  below the main diagonal containing nonzero elements, and  $n_u$  the number of diagonals of  $B$  above the main diagonal containing nonzero elements. The argument  $ka$  is the number of rows in the dimension statement for  $A$  in the calling program, and  $kb$  the number of rows in the dimension statement for  $B$  in the calling program. It is assumed that  $ka \geq m$  and  $kb \geq n$ .

It is assumed that  $AB$  is to be stored in band form in  $C$ .  $C$  is a 2-dimensional array of dimension  $kc \times nc$  where  $kc \geq m$ . The input argument  $nc$  is an estimate of the maximum number of diagonals of  $AB$  which will have to be stored ( $nc \leq \min\{n-1, m_\ell+n_\ell\} + \min\{\ell-1, m_u+n_u\} + 1$ ). BPROD is used if  $A, B$ , and  $C$  are real arrays, and CBPROD is used if  $A, B$ , and  $C$  are complex arrays. IERR,  $\nu_\ell$ , and  $\nu_u$  are integer variables. When the routine is called, if  $nc$  specifies sufficient storage for  $C$  then  $AB$  is computed and stored in band form in  $C$ . Also, IERR is assigned the value 0,  $\nu_\ell$  = the number of diagonals of  $AB$  below the main diagonal containing nonzero elements, and  $\nu_u$  = the number of diagonals of  $AB$  above the main diagonal containing nonzero elements.

**Error Return.** If  $nc$  does not specify sufficient storage for  $C$ , then IERR is assigned the value  $\nu$  where  $\nu$  is an estimate of the number of columns needed for  $C$ . Reset  $nc \geq \nu$ .

**Note.** It is assumed that the storage area for  $C$  does not overlap with the storage areas for  $A$  and  $B$ .

**Programmer.** A. H. Morris

## PRODUCT OF A REAL BANDED MATRIX AND VECTOR

Let  $A$  be a real  $m \times n$  matrix stored in band form. Then the following subroutines are available for multiplying  $A$  with a real vector.

**CALL BVPRD**( $m, n, A, ka, m_l, m_u, x, y$ )

$A$  is an  $m \times n$  matrix stored in band form,  $m_l$  the number of diagonals below the main diagonal containing nonzero elements, and  $m_u$  the number of diagonals above the main diagonal containing nonzero elements. The argument  $ka$  is the number of rows in the dimension statement for  $A$  in the calling program. It is assumed that  $0 \leq m_l < m, 0 \leq m_u < n$ , and  $ka \geq m$ .

The argument  $x$  is a column vector of dimension  $n$  and  $y$  an array of dimension  $m$ . When BVPRD is called, the product  $Ax$  is computed and stored in  $y$ .

**Remark.** It is assumed that the arrays  $A, x, y$  do not overlap.

**Programmer.** A. H. Morris

**CALL BVPRD1**( $m, n, A, ka, m_l, m_u, x, y$ )

$A$  is an  $m \times n$  matrix stored in band form,  $m_l$  the number of diagonals below the main diagonal containing nonzero elements, and  $m_u$  the number of diagonals above the main diagonal containing nonzero elements. The argument  $ka$  is the number of rows in the dimension statement for  $A$  in the calling program. It is assumed that  $0 \leq m_l < m, 0 \leq m_u < n$ , and  $ka \geq m$ .

The arguments  $x$  and  $y$  are column vectors of dimension  $n$  and  $m$  respectively. When BVPRD1 is called,  $Ax + y$  is computed and stored in  $y$ .

**Remark.** It is assumed that the arrays  $A, x, y$  do not overlap.

**Programmer.** A. H. Morris

**CALL BTPRD**( $m, n, A, ka, m_l, m_u, x, y$ )

$A$  is an  $m \times n$  matrix stored in band form,  $m_l$  the number of diagonals below the main diagonal containing nonzero elements, and  $m_u$  the number of diagonals above the main diagonal containing nonzero elements. The argument  $ka$  is the number of rows in the dimension statement for  $A$  in the calling program. It is assumed that  $0 \leq m_l < m, 0 \leq m_u < n$ , and  $ka \geq m$ .

The argument  $x$  is a row vector of dimension  $m$  and  $y$  an array of dimension  $n$ . When BTPRD is called, the product  $xA$  is computed and stored in  $y$ .

**Remark.** It is assumed that the arrays  $A, x, y$  do not overlap.



**Programmer.** A. H. Morris

**CALL** BTPRD1( $m, n, A, ka, m_l, m_u, x, y$ )

$A$  is an  $m \times n$  matrix stored in band form,  $m_l$  the number of diagonals below the main diagonal containing nonzero elements, and  $m_u$  then number of diagonals above the main diagonal containing nonzero elements. The argument  $ka$  is the number of rows in the dimension statement for  $A$  in the calling program. It is assumed that  $0 \leq m_l < m, 0 \leq m_u < n$ , and  $ka \geq m$ .

The arguments  $x$  and  $y$  are row vectors of dimension  $m$  and  $n$  respectively. When BTPRD1 is called,  $xA + y$  is computed and stored in  $y$ .

**Remark.** It is assumed that the arrays  $A, x, y$  do not overlap.

**Programmer.** A. H. Morris

## PRODUCT OF A COMPLEX BANDED MATRIX AND VECTOR

Let  $A$  be a complex  $m \times n$  matrix stored in band form. Then the following subroutines are available for multiplying  $A$  with a complex vector.

**CALL CBVPD**( $m, n, A, ka, m_l, m_u, x, y$ )

$A$  is an  $m \times n$  matrix stored in band form,  $m_l$  the number of diagonals below the main diagonal containing nonzero elements, and  $m_u$  the number of diagonals above the main diagonal containing nonzero elements. The argument  $ka$  is the number of rows in the dimension statement for  $A$  in the calling program. It is assumed that  $0 \leq m_l < m, 0 \leq m_u < n$ , and  $ka \geq m$ .

The argument  $x$  is a column vector of dimension  $n$  and  $y$  an array of dimension  $m$ .  $A, x, y$  are complex arrays. When CBVPD is called,  $Ax$  is computed and stored in  $y$ .

**Remark.** It is assumed that the arrays  $A, x, y$  do not overlap.

**Programmer.** A. H. Morris

**CALL CBVPD1**( $m, n, A, ka, m_l, m_u, x, y$ )

$A$  is an  $m \times n$  matrix stored in band form,  $m_l$  the number of diagonals below the main diagonal containing nonzero elements, and  $m_u$  the number of diagonals above the main diagonal containing nonzero elements. The argument  $ka$  is the number of rows in the dimension statement for  $A$  in the calling program. It is assumed that  $0 \leq m_l < m, 0 \leq m_u < n$ , and  $ka \geq m$ .

The arguments  $x$  and  $y$  are column vectors of dimension  $n$  and  $m$  respectively.  $A, x, y$  are complex arrays. When CBVPD1 is called,  $Ax + y$  is computed and stored in  $y$ .

**Remark.** It is assumed that the arrays  $A, x, y$  do not overlap.

**Programmer.** A. H. Morris

**CALL CBTPD**( $m, n, A, ka, m_l, m_u, x, y$ )

$A$  is an  $m \times n$  matrix stored in band form,  $m_l$  the number of diagonals below the main diagonal containing nonzero elements, and  $m_u$  the number of diagonals above the main diagonal containing nonzero elements. The argument  $ka$  is the number of rows in the dimension statement for  $A$  in the calling program. It is assumed that  $0 \leq m_l < m, 0 \leq m_u < n$ , and  $ka \geq m$ .

The argument  $x$  is a row vector of dimension  $m$  and  $y$  an array of dimension  $n$ .  $A, x, y$  are complex arrays. When CBTPD is called,  $xA$  is computed and stored in  $y$ .

**Remark.** It is assumed that the arrays  $A, x, y$  do not overlap.

Programmer. A. H. Morris

**CALL CBTPD1**( $m, n, A, ka, m_\ell, m_u, x, y$ )

$A$  is an  $m \times n$  matrix stored in band form,  $m_\ell$  the number of diagonals below the main diagonal containing nonzero elements, and  $m_u$  the number of diagonals above the main diagonal containing nonzero elements. The argument  $ka$  is the number of rows in the dimension statement for  $A$  in the calling program. It is assumed that  $0 \leq m_\ell < m, 0 \leq m_u < n$ , and  $ka \geq m$ .

The arguments  $x$  and  $y$  are row vectors of dimension  $m$  and  $n$  respectively.  $A, x, y$  are complex arrays. When CBTPD1 is called,  $xA + y$  is computed and stored in  $y$ .

**Remark.** It is assumed that the arrays  $A, x, y$  do not overlap.

Programmer. A. H. Morris

## SOLUTION OF BANDED SYSTEMS OF REAL LINEAR EQUATIONS

Let  $A$  be a nonsingular  $n \times n$  real matrix stored in band form and  $b$  a real column vector of dimension  $n$ . The subroutine BSLV is available for solving the system of equations  $Ax = b$ , and the subroutine BSLV1 is available for solving the transposed system of equations  $A^t x = b$ . On an initial call to either routine, partial pivot Gauss elimination is first employed to obtain an LU decomposition of  $A$ , and then the equations are solved. BSLV and BSLV1 always generate the same LU decomposition of  $A$ . After the decomposition is obtained on an initial call to BSLV or BSLV1, either routine may be called to solve a new system of equations  $Ax = r$  or  $A^t x = r$  without having to redecompose the matrix  $A$ .

```
CALL BSLV(MO,A,ka,n,ml,mu,X,IWK,IND)
CALL BSLV1(MO,A,ka,n,ml,mu,X,IWK,IND)
```

BSLV is called for solving  $Ax = b$  and BSLV1 is called for solving  $A^t x = b$ . The argument  $m_l$  is the number of diagonals below the main diagonal of  $A$  containing nonzero elements, and  $m_u$  the number of diagonals above the main diagonal containing nonzero elements. It is assumed that  $n \geq 1$ ,  $0 \leq m_l < n$ , and  $0 \leq m_u < n$ . MO is an input argument which specifies if BSLV or BSLV1 is being called for the first time. On an initial call, MO = 0 and we have the following setup:

$A$  is a 2-dimensional array of dimension  $ka \times m$  where  $ka \geq n$  and  $m \geq 2m_l + m_u + 1$ . On input, the first  $m_l + m_u + 1$  columns of the array contain the matrix  $A$  in band form. When the routine terminates, the array  $A$  will contain the upper triangular matrix  $U$  of the LU decomposition and the multipliers which were used to obtain it.

$X$  is an array of dimension  $n$  or larger. On input,  $X$  contains the vector  $b$ . On output,  $X$  will contain the solution of the system of equations.

IWK is an array of dimension  $n$  or larger for internal use by the routine. The pivot indices involved in the LU decomposition are stored in IWK.

On an initial call to BSLV or BSLV1, IND is an integer variable that reports the status of the results. When the routine terminates, IND has one of the following values:

```
IND = 0 The system of equations was solved.
IND = -1 Either  $n \leq 0$  or  $ka < n$ .
IND = -2 Either  $m_l < 0$  or  $m_l \geq n$ .
IND = -3 Either  $m_u < 0$  or  $m_u \geq n$ .
IND = k Column  $k$  of  $A$  has been reduced to a column containing only
        zeros.
```

After an initial call to BSLV or BSLV1, if IND = 0 on output then either routine may be called with MO  $\neq$  0. When MO  $\neq$  0 it is assumed that only  $b$  may have been modified. BSLV is called for solving the new set of equations  $Ax = b$ , and BSLV1 is called for solving the new set of equations  $A^t x = b$ . The routine employs the LU decomposition obtained on the initial call to BSLV or BSLV1 to solve the new set of equations. On input,  $X$  contains

the new vector  $b$ . On output,  $X$  will contain the solution to the new set of equations. In this case, IND is not referenced by the routine.

**Programming.** BSLV and BSLV1 employ the subroutines SNBFA, SNBSL, SAXPY, SSCAL, SSWAP and the functions ISAMAX and SDOT. SNBFA and SNBSL were written by E. A. Voorhees (Los Alamos Scientific Laboratory) and modified by A. H. Morris. The original versions of SNBFA and SNBSL are distributed by the SLATEC library.

## SOLUTION OF BANDED SYSTEMS OF COMPLEX LINEAR EQUATIONS

Let  $A$  be a nonsingular  $n \times n$  complex matrix stored in band form and  $b$  a complex column vector of dimension  $n$ . The subroutine CBSLV is available for solving the system of equations  $Ax = b$ , and the subroutine CBSLV1 is available for solving the transposed system of equations  $A^t x = b$ . On an initial call to either routine, partial pivot Gauss elimination is first employed to obtain an LU decomposition of  $A$ , and then the equations are solved. CBSLV and CBSLV1 always generate the same LU decomposition of  $A$ . After the decomposition is obtained on an initial call to CBSLV or CBSLV1, either routine may be called to solve a new system of equations  $Ax = r$  or  $A^t x = r$  without having to redecompose the matrix  $A$ .

```
CALL CBSLV(MO,A,ka,n,ml,mu,X,IWK,IND)
CALL CBSLV1(MO,A,ka,n,ml,mu,X,IWK,IND)
```

CBSLV is called for solving  $Ax = b$  and CBSLV1 is called for solving  $A^t x = b$ . The argument  $m_l$  is the number of diagonals below the main diagonal of  $A$  containing nonzero elements, and  $m_u$  the number of diagonals above the main diagonal containing nonzero elements. It is assumed that  $n \geq 1$ ,  $0 \leq m_l < n$ , and  $0 \leq m_u < n$ . MO is an input argument which specifies if CBSLV or CBSLV1 is being called for the first time. On an initial call, MO = 0 and we have the following setup:

$A$  is a 2-dimensional array of dimension  $ka \times m$  where  $ka \geq n$  and  $m \geq 2m_l + m_u + 1$ . On input, the first  $m_l + m_u + 1$  columns of the array contain the matrix  $A$  in band form. When the routine terminates, the array  $A$  will contain the upper triangular matrix  $U$  of the LU decomposition and the multipliers which were used to obtain it.

$X$  is an array of dimension  $n$  or larger. On input,  $X$  contains the vector  $b$ . On output,  $X$  will contain the solution of the system of equations.

IWK is an array of dimension  $n$  or larger for internal use by the routine. The pivot indices involved in the LU decomposition are stored in IWK.

On an initial call to CBSLV or CBSLV1, IND is an integer variable that reports the status of the results. When the routine terminates, IND has one of the following values:

```
IND = 0 The system of equations was solved.
IND = -1 Either  $n \leq 0$  or  $ka < n$ .
IND = -2 Either  $m_l < 0$  or  $m_l \geq n$ .
IND = -3 Either  $m_u < 0$  or  $m_u \geq n$ .
IND = k Column  $k$  of  $A$  has been reduced to a column containing only
        zeros.
```

After an initial call to CBSLV or CBSLV1, if IND = 0 on output then either routine may be called with MO  $\neq$  0. When MO  $\neq$  0 it is assumed that only  $b$  may have been modified. CBSLV is called for solving the new set of equations  $Ax = b$ , and CBSLV1 is called for solving the new set of equations  $A^t x = b$ . The routine employs the LU decomposition

obtained on the initial call to CBSLV or CBSLV1 to solve the new set of equations. On input,  $X$  contains the new vector  $b$ . On output,  $X$  will contain the solution to the new set of equations. In this case, IND is not referenced by the routine.

**Programming.** CBSLV and CBSLV1 employ the subroutines CBFA, CBSL, CAXPY, CSCAL, CSWAP and the functions ICAMAX and CDOTU. CBFA and CBSL are adaptations by A. H. Morris of the subroutines SNBFA and SNBSL, written by E. A. Voorhees (Los Alamos Scientific Laboratory). SNBFA and SNBSL are distributed by the SLATEC library.

## STORAGE OF SPARSE MATRICES

A matrix is said to be *sparse* if it contains sufficiently many zero elements for it to be worthwhile to use special techniques that avoid storing and operating with the zeros. The scheme adopted by the NSWCL library for storing a sparse  $m \times n$  matrix  $(a_{ij})$  requires three 1-dimensional arrays  $A$ ,  $IA$ ,  $JA$ . The array  $A$  contains the nonzero elements of the matrix, stored row by row. The array  $JA$  contains the column numbers of the corresponding elements of the  $A$  array; i.e., if  $A(k)$  contains  $(a_{ij})$  then  $JA(k) = j$ . The elements of a row of the matrix may be given in any order in  $A$ .

$IA$  is an array containing  $m+1$  integers which specify where the rows of the matrix are stored in  $A$ . For  $i \leq m$ ,  $IA(i)$  is the index of the location in  $A$  where the  $i^{th}$  row information begins. It is assumed that the rows are stored sequentially; i.e., that  $IA(1) \leq \dots \leq IA(m)$ .  $IA(m+1)$  is set so that  $IA(m+1) - IA(1) =$  the number of elements stored in  $A$ . For  $i \leq m$ , if  $IA(i) < IA(i+1)$  then  $A(\ell)$  is the first entry of the  $i^{th}$  row of the matrix in  $A$  where  $\ell = IA(i)$ . Otherwise, if  $IA(i) = IA(i+1)$  then no entries for the  $i^{th}$  row of the matrix are stored in  $A$ . This can occur only if the  $i^{th}$  row of the matrix consists entirely of zeros. If this occurs then the  $i^{th}$  row is called a **null row** of  $A$ . For any  $i \leq m$ ,  $IA(i+1) - IA(i)$  is the number of entries for the  $i^{th}$  row of the matrix that are stored in  $A$ . For convenience,  $IA(i+1) - IA(i)$  is called the **length** of the  $i^{th}$  row.

**Example.** The matrix

$$\begin{pmatrix} a_{11} & a_{12} & 0 & 0 & 0 & 0 & 0 & a_{18} \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & a_{37} & a_{38} \\ 0 & 0 & a_{43} & 0 & 0 & 0 & 0 & 0 \end{pmatrix}$$

can be stored as follows:

$A:$	<table><tr><td><math>a_{11}</math></td><td><math>a_{18}</math></td><td><math>a_{12}</math></td><td><math>a_{37}</math></td><td><math>a_{38}</math></td><td><math>a_{43}</math></td></tr></table>	$a_{11}$	$a_{18}$	$a_{12}$	$a_{37}$	$a_{38}$	$a_{43}$
$a_{11}$	$a_{18}$	$a_{12}$	$a_{37}$	$a_{38}$	$a_{43}$		
$JA:$	<table><tr><td>1</td><td>8</td><td>2</td><td>7</td><td>8</td><td>3</td></tr></table>	1	8	2	7	8	3
1	8	2	7	8	3		
$IA:$	<table><tr><td>1</td><td>4</td><td>4</td><td>6</td><td>7</td></tr></table>	1	4	4	6	7	
1	4	4	6	7			

The storage of the elements  $a_{11}$   $a_{12}$   $a_{18}$  in the order  $a_{11}$   $a_{18}$   $a_{12}$  is permissible. The elements of a row of the matrix may be given in any order desired.

**Remark.** It is not required that each  $a_{ij}$  in  $A$  be nonzero.



## CONVERSION OF SPARSE MATRICES TO AND FROM THE STANDARD FORMAT

The following subroutines permit one to convert sparse matrices to and from the standard format.

```
CALL CVRS(A, ka, m, n, B, IB, JB, NUM, IERR)
CALL CVDS(A, ka, m, n, B, IB, JB, NUM, IERR)
CALL CVCS(A, ka, m, n, B, IB, JB, NUM, IERR)
```

$A$  is an  $m \times n$  matrix stored in the standard format. The argument  $ka$  is the number of rows in the dimension statement for  $A$  in the calling program. It is assumed that  $A$  is to be stored in sparse form in the arrays  $B$ ,  $IB$ ,  $JB$ . CVRS is used if  $A$  is a real matrix and  $B$  a real array, CVDS is used if  $A$  is a double precision matrix and  $B$  a double precision array, and CVCS is used if  $A$  is a complex matrix and  $B$  a complex array.

The input argument NUM is the estimated maximum number of elements that will appear in  $B$  and  $JB$ . It is assumed that  $B$  and  $JB$  are of dimension  $\max\{1, \text{NUM}\}$  and that  $IB$  is of dimension  $m + 1$ . IERR is an integer variable. When the routine is called, if NUM specifies sufficient storage for  $B$  and  $JB$ , then  $A$  is stored in  $B$ ,  $IB$ ,  $JB$  and IERR is assigned the value 0.

**Error Return.** If it is found that there is not sufficient storage in  $B$  and  $JB$  for the  $i^{\text{th}}$  row of  $A$ , then IERR is set to  $i$  and the routine terminates. In this case, if  $i > 1$  then the first  $i - 1$  rows of  $A$  will have been stored in  $B$  and  $JB$ , and  $IB(1), \dots, IB(i)$  will contain the appropriate row locations.

**Remark.** No zero elements of  $A$  are stored in  $B$ , and the elements of each row of  $B$  are ordered so that the column indices of the elements of the row are in ascending order.

**Example.** If  $A$  is the  $m \times n$  zero matrix then NUM can be set to 0. In this case the result will be  $IB(1) = \dots = IB(m + 1) = 1$ .

**Note.** It is assumed that the storage areas  $A$  and  $B$  do not overlap.

**Programmer.** A. H. Morris.

```
CALL CVSR(A, IA, JA, B, kb, m, n)
CALL CVSD(A, IA, JA, B, kb, m, n)
CALL CVSC(A, IA, JA, B, kb, m, n)
```

$A$  is an  $m \times n$  sparse matrix stored in the arrays  $A$ ,  $IA$ ,  $JA$ , and  $B$  is a 2-dimensional array of dimension  $kb \times n$  where  $kb \geq m$ . CVSR is used if  $A$  and  $B$  are real arrays, CVSD is used if  $A$  and  $B$  are double precision arrays, and CVSC is used if  $A$  and  $B$  are complex arrays. When the routine is called, the matrix  $A$  is stored in the array  $B$  in the standard format.

**Note.** It is assumed that the storage areas  $A$  and  $B$  do not overlap.

**Programmer.** A. H. Morris.

## CONVERSION OF SPARSE REAL MATRICES TO AND FROM DOUBLE PRECISION FORM

The following subroutines are available for converting sparse real matrices to and from double precision form.

**CALL SCVRD(*A, IA, JA, B, IB, JB, m*)**

*A* is a sparse real matrix stored in the arrays *A, IA, JA*. *A* is a real array and *B* a double precision array. If *A* and *JA* contain *k* elements then it is assumed that *B* and *JB* are arrays of dimension *k*. It is also assumed that the matrix *A* has  $m \geq 1$  rows and that *IB* is an array of dimension  $m + 1$ . SCVRD stores the matrix *A* in double precision form in *B, IB, JB*.

### Remarks.

- (1) No zero elements of *A* are stored in *B*.
- (2) It is assumed that the arrays *IA, JA* and *IB, JB* reference different storage areas.

Programmer. A. H. Morris.

**CALL SCVDR(*A, IA, JA, B, IB, JB, m*)**

*A* is a sparse double precision matrix stored in the arrays *A, IA, JA*. *A* is a double precision array and *B* a real array. If *A* and *JA* contain *k* elements then it is assumed that *B* and *JB* are arrays of dimension *k*. It is also assumed that the matrix *A* has  $m \geq 1$  rows and that *IB* is an array of dimension  $m + 1$ . SCVDR stores the matrix *A* in single precision form in *B, IB, JB*.

### Remarks.

- (1) No zero elements of *A* are stored in *B*.
- (2) It is assumed that the arrays *IA, JA* and *IB, JB* reference different storage areas.

Programmer. A. H. Morris.

## THE REAL AND IMAGINARY PARTS OF A SPARSE COMPLEX MATRIX

If  $A = (A_{ij})$  is a complex matrix then let  $\text{Re}(A) = (\text{Re}(a_{ij}))$  and  $\text{Im}(A) = (\text{Im}(a_{ij}))$  denote the real and imaginary parts of  $A$ . If the matrix  $A$  is stored in sparse form, then the following subroutines are available for obtaining  $\text{Re}(A)$  and  $\text{Im}(A)$  in sparse form.

**CALL CSREAL( $A, IA, JA, B, IB, JB, m$ )**

$A$  is a sparse complex matrix stored in the arrays  $A, IA, JA$ .  $A$  is a complex array and  $B$  a real array. If  $A$  and  $JA$  contain  $k$  elements then it is assumed that  $B$  and  $JB$  are arrays of dimension  $k$ . It is also assumed that the matrix  $A$  has  $m \geq 1$  rows and that  $IB$  is an array of dimension  $m + 1$ . CSREAL stores  $\text{Re}(A)$  in sparse form in  $B, IB, JB$ .

### Remarks.

- (1) No zero elements of  $\text{Re}(A)$  are stored in  $B$ .
- (2) It is assumed that the arrays  $IA, JA$  and  $IB, JB$  reference different storage areas.

**Programmer.** A. H. Morris.

**CALL CSIMAG( $A, IA, JA, B, IB, JB, m$ )**

$A$  is a sparse complex matrix stored in the arrays  $A, IA, JA$ .  $A$  is a complex array and  $B$  a real array. If  $A$  and  $JA$  contain  $k$  elements then it is assumed that  $B$  and  $JB$  are arrays of dimension  $k$ . It is also assumed that the matrix  $A$  has  $m \geq 1$  rows and that  $IB$  is an array of dimension  $m + 1$ . CSIMAG stores  $\text{Im}(A)$  in sparse form in  $B, IB, JB$ .

### Remarks.

- (1) No zero elements of  $\text{Im}(A)$  are stored in  $B$ .
- (2) It is assumed that the arrays  $IA, JA$  and  $IB, JB$  reference different storage areas.

**Programmer.** A. H. Morris.

## COMPUTING $A + Bi$ FOR SPARSE REAL MATRICES $A$ AND $B$

Given the real  $m \times n$  matrices  $A$  and  $B$  stored in sparse form. Then the subroutine SCVRC is available for obtaining the complex matrix  $A + Bi$  where  $i = \sqrt{-1}$ .

**CALL SCVRC**( $A, IA, JA, B, IB, JB, C, IC, JC, m, n, NUM, WK, IERR$ )

$A$  and  $B$  are real  $m \times n$  matrices stored in the arrays  $A, IA, JA$  and  $B, IB, JB$ . It is assumed that  $A + Bi$  is to be stored in sparse form in the arrays  $C, IC, JC$ .  $A$  and  $B$  are real arrays and  $C$  a complex array.  $NUM$  is the estimated maximum number of elements that will appear in  $C$  and  $JC$ . It is assumed that  $C$  and  $JC$  are arrays of dimension  $\max\{1, NUM\}$  and that  $IC$  is an array of dimension  $m + 1$ .

$WK$  is a real array of dimension  $n$  or larger that is a work space for the routine.

$IERR$  is an integer variable. When SCVRC is called, if  $NUM$  specifies sufficient storage for  $C$  and  $JC$  then  $A + Bi$  is stored in  $C, IC, JC$ . Also  $IERR$  is assigned the value 0.

**Error Return.** If there is not sufficient storage in  $C$  and  $JC$  for the  $k^{th}$  row of  $A + Bi$ , then  $IERR$  is set to  $k$  and the routine terminates. In this case, if  $k > 1$  then the first  $k - 1$  rows of  $A + Bi$  will have been stored in  $C$  and  $JC$ . Also  $IC(1), \dots, IC(k)$  will contain the appropriate row locations.

**Remark.** No zeros are stored in  $C$ .

**Programmer.** A. H. Morris

## COPYING SPARSE MATRICES

The following subroutines are available for copying sparse matrices.

```
CALL RSCOPY(A,IA,JA,B,IB,JB,m)
CALL DSCOPY(A,IA,JA,B,IB,JB,m)
CALL CSCOPY(A,IA,JA,B,IB,JB,m)
```

RSCOPY is used if  $A$  and  $B$  are real arrays, DSCOPY is used if  $A$  and  $B$  are double precision arrays, and CSCOPY is used if  $A$  and  $B$  are complex arrays.

$A$  is a sparse matrix stored in the arrays  $A, IA, JA$ . If  $A$  and  $JA$  contain  $k$  elements then it is assumed that  $B$  and  $JB$  are arrays of dimension  $k$ . It is also assumed that the matrix  $A$  has  $m \geq 1$  rows and that  $IB$  is an array of dimension  $m + 1$ . The routine copies the matrix  $A$  and stores the copy in  $B, IB, JB$ .

### Remarks.

- (1) No zero elements of  $A$  are stored in  $B$ .
- (2) It is assumed that the arrays  $A, IA, JA$  and  $B, IB, JB$  reference different storage areas.

Programmer. A. H. Morris.

## COMPUTING CONJUGATES OF SPARSE COMPLEX MATRICES

If  $A = (a_{ij})$  is a complex matrix stored in sparse form, then the following subroutine is available for computing the conjugate matrix  $\bar{A} = (\bar{a}_{ij})$ .

**CALL SCONJ**( $A, IA, JA, B, IB, JB, m$ )

It is assumed that the sparse complex matrix  $A$  is stored in the arrays  $A, IA, JA$ . If  $A$  and  $JA$  contain  $k$  elements, then it is also assumed that  $B$  and  $JB$  are arrays of dimension  $k$ .  $A$  and  $B$  are complex arrays. It is assumed that the matrix  $A$  has  $m \geq 1$  rows and that  $IB$  is an array of dimension  $m + 1$ . When the routine is called, the conjugate matrix  $\bar{A}$  is stored in  $B, IB, JB$ .

**Remark.** The user may let  $B = A, IB = IA$ , and  $JB = JA$  when  $IA(1) = 1$ .

**Programmer.** A. H. Morris

## TRANSPOSING SPARSE REAL MATRICES

The subroutines RPOSE and RPOSE1 are available for transposing a sparse  $m \times n$  real matrix  $A$ . RPOSE1 is more general than RPOSE. For any permutation  $\pi = \{i_1, \dots, i_m\}$  of  $\{1, \dots, m\}$  let  $P$  denote the corresponding  $m \times m$  permutation matrix. Then RPOSE1 computes the matrix  $(PA)^t$ .

**CALL RPOSE**( $A, IA, JA, B, IB, JB, m, n$ )

It is assumed that the sparse matrix  $A$  is stored in the arrays  $A, IA, JA$ . If  $A$  and  $JA$  contain  $k$  elements, then it is also assumed that  $B$  and  $JB$  are arrays of dimension  $k$  and that  $IB$  is an array of dimension  $n + 1$ . When RPOSE is called, the transpose  $A^t$  is stored in  $B, IB, JB$ .

**Remarks.** RPOSE orders the elements of each row of  $A^t$  so that the column indices of the elements of the row are in ascending order. However, it does no checking for zero elements in  $A$ . If zero elements appear in the array  $A$ , then the zero elements will also appear in  $B$ .

**Restriction.** It is assumed that the storage areas  $B, IB, JB$  do not overlap with the storage areas  $A, IA, JA$ .

**Programmer.** A. H. Morris

**Reference.** Gustavson, F. G., "Two Fast Algorithms for Sparse Matrices: Multiplication and Permuted Transposition," *ACM Trans. Math Software* 4 (1978), pp. 250-269.

**CALL RPOSE1**( $\pi, A, IA, JA, B, IB, JB, m, n$ )

It is assumed that  $\pi$  is an integer array of dimension  $m$  containing the data  $\{i_1, \dots, i_m\}$ , and that the sparse matrix  $A$  is stored in the arrays  $A, IA, JA$ . If  $A$  and  $JA$  contain  $k$  elements, then it is also assumed that  $B$  and  $JB$  are arrays of dimension  $k$  and that  $IB$  is an array of dimension  $n + 1$ . When RPOSE1 is called,  $(PA)^t$  is computed and the results are stored in  $B, IB, JB$ .

**Remarks.** RPOSE1 orders the elements of each row of  $(PA)^t$  so that the column indices of the elements of the row are in ascending order. However, it does no checking for zero elements in  $A$ . If zero elements appear in the array  $A$ , then the zero elements will also appear in  $B$ .

**Restriction.** It is assumed that the storage areas  $B, IB, JB$  do not overlap with the storage areas  $A, IA, JA$ .

**Programmer.** A. H. Morris

**Reference.** Gustavson, F. G., "Two Fast Algorithms for Sparse Matrices: Multiplication and Permuted Transposition," *ACM Trans. Math Software* 4 (1978), pp. 250-269.



## TRANSPOSING SPARSE DOUBLE PRECISION MATRICES

The subroutines DPOSE and DPOSE1 are available for transposing a sparse  $m \times n$  double precision matrix  $A$ . DPOSE1 is more general than DPOSE. For any permutation  $\pi = \{i_1, \dots, i_m\}$  of  $\{1, \dots, m\}$  let  $P$  denote the corresponding  $m \times m$  permutation matrix. Then DPOSE1 computes the matrix  $(PA)^t$ .

**CALL DPOSE**( $A, IA, JA, B, IB, JB, m, n$ )

It is assumed that the sparse matrix  $A$  is stored in the arrays  $A, IA, JA$ .  $A$  and  $B$  are double precision arrays. If  $A$  and  $JA$  contain  $k$  elements, then it is also assumed that  $B$  and  $JB$  are arrays of dimension  $k$  and that  $IB$  is an array of dimension  $n + 1$ . When DPOSE is called, the transpose  $A^t$  is stored in  $B, IB, JB$ .

**Remarks.** DPOSE orders the elements of each row of  $A^t$  so that the column indices of the elements of the row are in ascending order. However, it does no checking for zero elements in  $A$ . If zero elements appear in the array  $A$ , then the zero elements will also appear in  $B$ .

**Restriction.** It is assumed that the storage areas  $B, IB, JB$  do not overlap with the storage areas  $A, IA, JA$ .

**Programmer.** A. H. Morris

**Reference.** Gustavson, F. G., "Two Fast Algorithms for Sparse Matrices: Multiplication and Permuted Transposition," *ACM Trans. Math Software* 4 (1978), pp. 250-269.

**CALL DPOSE1**( $\pi, A, IA, JA, B, IB, JB, m, n$ )

It is assumed that  $\pi$  is an integer array of dimension  $m$  containing the data  $\{i_1, \dots, i_m\}$ , and that the sparse matrix  $A$  is stored in the arrays  $A, IA, JA$ .  $A$  and  $B$  are double precision arrays. If  $A$  and  $JA$  contain  $k$  elements, then it is also assumed that  $B$  and  $JB$  are arrays of dimension  $k$  and that  $IB$  is an array of dimension  $n + 1$ . When DPOSE1 is called,  $(PA)^t$  is computed and the results are stored in  $B, IB, JB$ .

**Remarks.** DPOSE1 orders the elements of each row of  $(PA)^t$  so that the column indices of the elements of the row are in ascending order. However, it does no checking for zero elements in  $A$ . If zero elements appear in the array  $A$ , then the zero elements will also appear in  $B$ .

**Restriction.** It is assumed that the storage areas  $B, IB, JB$  do not overlap with the storage areas  $A, IA, JA$ .

**Programmer.** A. H. Morris

**Reference.** Gustavson, F. G., "Two Fast Algorithms for Sparse Matrices: Multiplication and Permuted Transposition," *ACM Trans. Math Software* 4 (1978), pp. 250-269.

## TRANSPOSING SPARSE COMPLEX MATRICES

The subroutines CPOSE and CPOSE1 are available for transposing a sparse  $m \times n$  complex matrix  $A$ . CPOSE1 is more general than CPOSE. For any permutation  $\pi = \{i_1, \dots, i_m\}$  of  $\{1, \dots, m\}$  let  $P$  denote the corresponding  $m \times m$  permutation matrix. Then CPOSE1 computes the matrix  $(PA)^t$ .

**CALL CPOSE**( $A, IA, JA, B, IB, JB, m, n$ )

It is assumed that the sparse matrix  $A$  is stored in the arrays  $A, IA, JA$ .  $A$  and  $B$  are complex arrays. If  $A$  and  $JA$  contain  $k$  elements, then it is also assumed that  $B$  and  $JB$  are arrays of dimension  $k$  and that  $IB$  is an array of dimension  $n + 1$ . When CPOSE is called, the transpose  $A^t$  is stored in  $B, IB, JB$ .

**Remarks.** CPOSE orders the elements of each row of  $A^t$  so that the column indices of the elements of the row are in ascending order. However, it does no checking for zero elements in  $A$ . If zero elements appear in the array  $A$ , then the zero elements will also appear in  $B$ .

**Restriction.** It is assumed that the storage areas  $B, IB, JB$  do not overlap with the storage areas  $A, IA, JA$ .

**Programmer.** A. H. Morris

**Reference.** Gustavson, F. G., "Two Fast Algorithms for Sparse Matrices: Multiplication and Permuted Transposition," *ACM Trans. Math Software* 4 (1978), pp. 250-269.

**CALL CPOSE1**( $\pi, A, IA, JA, B, IB, JB, m, n$ )

It is assumed that  $\pi$  is an integer array of dimension  $m$  containing the data  $\{i_1, \dots, i_m\}$ , and that the sparse matrix  $A$  is stored in the arrays  $A, IA, JA$ .  $A$  and  $B$  are complex arrays. If  $A$  and  $JA$  contain  $k$  elements, then it is also assumed that  $B$  and  $JB$  are arrays of dimension  $k$  and that  $IB$  is an array of dimension  $n + 1$ . When CPOSE1 is called,  $(PA)^t$  is computed and the results are stored in  $B, IB, JB$ .

**Remarks.** CPOSE1 orders the elements of each row of  $(PA)^t$  so that the column indices of the elements of the row are in ascending order. However, it does no checking for zero elements in  $A$ . If zero elements appear in the array  $A$ , then the zero elements will also appear in  $B$ .

**Restriction.** It is assumed that the storage areas  $B, IB, JB$  do not overlap with the storage areas  $A, IA, JA$ .

**Programmer.** A. H. Morris

**Reference.** Gustavson, F. G., "Two Fast Algorithms for Sparse Matrices: Multiplication and Permuted Transposition," *ACM Trans. Math Software* 4 (1978), pp. 250-269.

## ADDITION OF SPARSE MATRICES

The following subroutines are available for adding sparse matrices.

```
CALL SADD(A, IA, JA, B, IB, JB, C, IC, JC, m, n, NUM, WK, IERR)
CALL DSADD(A, IA, JA, B, IB, JB, C, IC, JC, m, n, NUM, WK, IERR)
CALL CSADD(A, IA, JA, B, IB, JB, C, IC, JC, m, n, NUM, WK, IERR)
```

$A$  and  $B$  are sparse  $m \times n$  matrices stored in the arrays  $A, IA, JA$  and  $B, IB, JB$ . It is assumed that  $A + B$  is to be stored in sparse form in the arrays  $C, IC, JC$ .  $NUM$  is the estimated maximum number of elements that will appear in  $C$  and  $JC$ . It is assumed that  $C$  and  $JC$  are arrays of dimension  $\max\{1, NUM\}$  and that  $IC$  is an array of dimension  $m + 1$ .

$SADD$  is used if  $A, B, C$ , and  $WK$  are real arrays,  $DSADD$  is used if  $A, B, C$ , and  $WK$  are double precision arrays, and  $CSADD$  is used if  $A, B, C$ , and  $WK$  are complex arrays.

$WK$  is an array of dimension  $n$  or larger that is a work space for the routine.

$IERR$  is an integer variable. When the routine is called, if  $NUM$  specifies sufficient storage for  $C$  and  $JC$  then  $A + B$  is computed and stored in  $C, IC, JC$ . Also  $IERR$  is assigned the value 0.

**Error Return.** If there is not sufficient storage in  $C$  and  $JC$  for the  $i^{th}$  row of  $A + B$ , then  $IERR$  is set to  $i$  and the routine terminates. In this case, if  $i > 1$  then the first  $i - 1$  rows of  $A + B$  will have been computed and stored in  $C$  and  $JC$ . Also  $IC(1), \dots, IC(i)$  will contain the appropriate row locations.

### Remarks.

- (1) No zeros are stored in  $C$ .
- (2) It is assumed that  $C, IC, JC$  reference different storage areas than  $A, IA, JA$  and  $B, IB, JB$ .

**Programmer.** A. H. Morris

## SUBTRACTION OF SPARSE MATRICES

The following subroutines are available for subtracting sparse matrices.

```
CALL SSUBT(A,IA,JA,B,IB,JB,C,IC,JC,m,n,NUM,WK,IERR)
CALL DSSUBT(A,IA,JA,B,IB,JB,C,IC,JC,m,n,NUM,WK,IERR)
CALL CSSUBT(A,IA,JA,B,IB,JB,C,IC,JC,m,n,NUM,WK,IERR)
```

$A$  and  $B$  are sparse  $m \times n$  matrices stored in the arrays  $A, IA, JA$  and  $B, IB, JB$ . It is assumed that  $A - B$  is to be stored in sparse form in the arrays  $C, IC, JC$ .  $NUM$  is the estimated maximum number of elements that will appear in  $C$  and  $JC$ . It is assumed that  $C$  and  $JC$  are arrays of dimension  $\max\{1, NUM\}$  and that  $IC$  is an array of dimension  $m + 1$ .

SSUBT is used if  $A, B, C$ , and  $WK$  are real arrays, DSSUBT is used if  $A, B, C$ , and  $WK$  are double precision arrays, and CSSUBT is used if  $A, B, C$ , and  $WK$  are complex arrays.

$WK$  is an array of dimension  $n$  or larger that is a work space for the routine.

$IERR$  is an integer variable. When the routine is called, if  $NUM$  specifies sufficient storage for  $C$  and  $JC$  then  $A - B$  is computed and stored in  $C, IC, JC$ . Also  $IERR$  is assigned the value 0.

**Error Return.** If there is not sufficient storage in  $C$  and  $JC$  for the  $i^{th}$  row of  $A - B$ , then  $IERR$  is set to  $i$  and the routine terminates. In this case, if  $i > 1$  then the first  $i - 1$  rows of  $A - B$  will have been computed and stored in  $C$  and  $JC$ . Also  $IC(1), \dots, C(i)$  will contain the appropriate row locations.

### Remarks.

- (1) No zeros are stored in  $C$ .
- (2) It is assumed that  $C, IC, JC$  reference different storage areas than  $A, IA, JA$  and  $B, IB, JB$ .

**Programmer.** A. H. Morris

## MULTIPLICATION OF SPARSE MATRICES

The following subroutines are available for multiplying sparse matrices.

```
CALL SPROD(A, IA, JA, B, IB, JB, C, IC, JC,  $\ell$ ,  $m$ ,  $n$ , NUM, WK, IERR)
CALL DSPROD(A, IA, JA, B, IB, JB, C, IC, JC,  $\ell$ ,  $m$ ,  $n$ , NUM, WK, IERR)
CALL CSPROD(A, IA, JA, B, IB, JB, C, IC, JC,  $\ell$ ,  $m$ ,  $n$ , NUM, WK, IERR)
```

$A$  is a sparse  $\ell \times m$  matrix stored in the arrays  $A, IA, JA$ , and  $B$  a sparse  $m \times n$  matrix stored in the arrays  $B, IB, JB$ . It is assumed that  $AB$  is to be stored in sparse form in the arrays  $C, IC, JC$ . NUM is the estimated maximum number of elements that will appear in  $C$  and  $JC$ . It is assumed that  $C$  and  $JC$  are arrays of dimension  $\max\{1, \text{NUM}\}$  and that  $IC$  is an array of dimension  $\ell + 1$ .

SPROD is used if  $A, B, C$ , and  $WK$  are real arrays, DSPROD is used if  $A, B, C$ , and  $WK$  are double precision arrays, and CSPROD is used if  $A, B, C$ , and  $WK$  are complex arrays.

$WK$  is an array of dimension  $n$  or larger that is a work space for the routine.

IERR is an integer variable. When the routine is called, if NUM specifies sufficient storage for  $C$  and  $JC$  then  $AB$  is computed and stored in  $C, IC, JC$ . Also IERR is assigned the value 0.

**Error Return.** If there is not sufficient storage in  $C$  and  $JC$  for the  $i^{\text{th}}$  row of  $AB$ , then IERR is set to  $i$  and the routine terminates. In this case, if  $i > 1$  then the first  $i - 1$  rows of  $AB$  will have been computed and stored in  $C$  and  $JC$ . Also  $IC(1), \dots, IC(i)$  will contain the appropriate row locations.

### Remarks.

- (1) No zeros are stored in  $C$ .
- (2) It is assumed that  $C, IC, JC$  reference different storage areas than  $A, IA, JA$  and  $B, IB, JB$ .

**Programmer.** A. H. Morris

## PRODUCT OF A REAL SPARSE MATRIX AND VECTOR

Let  $A$  be a real  $m \times n$  sparse matrix stored in the arrays  $A, IA, JA$ . Then the following subroutines are available for multiplying  $A$  with a real vector.

**CALL MVPRD( $m, n, A, IA, JA, x, y$ )**

The argument  $x$  is a column vector of dimension  $n$  and  $y$  an array of dimension  $m$ . When MVPRD is called,  $Ax$  is computed and stored in  $y$ .

**Remark.** It is assumed that the arrays  $A, x, y$  do not overlap.

**Programmer.** A. H. Morris

**CALL MVPRD1( $m, n, A, IA, JA, x, y$ )**

The arguments  $x$  and  $y$  are column vectors of dimension  $n$  and  $m$  respectively. When MVPRD1 is called,  $Ax + y$  is computed and stored in  $y$ .

**Remark.** It is assumed that the arrays  $A, x, y$  do not overlap.

**Programmer.** A. H. Morris

**CALL MTPRD( $m, n, A, IA, JA, x, y$ )**

The argument  $x$  is a row vector of dimension  $m$  and  $y$  an array of dimension  $n$ . When MTPRD is called,  $xA$  is computed and stored in  $y$ .

**Remark.** It is assumed that the arrays  $A, x, y$  do not overlap.

**Programmer.** A. H. Morris

**CALL MTPRD1( $m, n, A, IA, JA, x, y$ )**

The arguments  $x$  and  $y$  are row vectors of dimension  $m$  and  $n$  respectively. When MTPRD1 is called,  $xA + y$  is computed and stored in  $y$ .

**Remark.** It is assumed that the arrays  $A, x, y$  do not overlap.

**Programmer.** A. H. Morris

## PRODUCT OF A DOUBLE PRECISION SPARSE MATRIX AND VECTOR

Let  $A$  be a double precision  $m \times n$  sparse matrix stored in the arrays  $A, IA, JA$ . Then the following subroutines are available for multiplying  $A$  with a double precision vector.

**CALL DVPRD( $m, n, A, IA, JA, x, y$ )**

The argument  $x$  is a column vector of dimension  $n$  and  $y$  an array of dimension  $m$ .  $A, x, y$  are double precision arrays. When DVPRD is called,  $Ax$  is computed and stored in  $y$ .

**Remark.** It is assumed that the arrays  $A, x, y$  do not overlap.

**Programmer.** A. H. Morris

**CALL DVPRD1( $m, n, A, IA, JA, x, y$ )**

The arguments  $x$  and  $y$  are column vectors of dimension  $n$  and  $m$  respectively.  $A, x, y$  are double precision arrays. When DVPRD1 is called,  $Ax + y$  is computed and stored in  $y$ .

**Remark.** It is assumed that the arrays  $A, x, y$  do not overlap.

**Programmer.** A. H. Morris

**CALL DTPRD( $m, n, A, IA, JA, x, y$ )**

The argument  $x$  is a row vector of dimension  $m$  and  $y$  an array of dimension  $n$ .  $A, x, y$  are double precision arrays. When DTPRD is called,  $xA$  is computed and stored in  $y$ .

**Remark.** It is assumed that the arrays  $A, x, y$  do not overlap.

**Programmer.** A. H. Morris

**CALL DTPRD1( $m, n, A, IA, JA, x, y$ )**

The arguments  $x$  and  $y$  are row vectors of dimension  $m$  and  $n$  respectively.  $A, x, y$  are double precision arrays. When DTPRD1 is called,  $xA + y$  is computed and stored in  $y$ .

**Remark.** It is assumed that the arrays  $A, x, y$  do not overlap.

**Programmer.** A. H. Morris

## PRODUCT OF A COMPLEX SPARSE MATRIX AND VECTOR

Let  $A$  be a complex  $m \times n$  sparse matrix stored in the arrays  $A, IA, JA$ . Then the following subroutines are available for multiplying  $A$  with a complex vector.

**CALL CVPRD( $m, n, A, IA, JA, x, y$ )**

The argument  $x$  is a column vector of dimension  $n$  and  $y$  an array of dimension  $m$ .  $A, x, y$  are complex arrays. When CVPRD is called,  $Ax$  is computed and stored in  $y$ .

**Remark.** It is assumed that the arrays  $A, x, y$  do not overlap.

**Programmer.** A. H. Morris

**CALL CVPRD1( $m, n, A, IA, JA, x, y$ )**

The arguments  $x$  and  $y$  are column vectors of dimension  $n$  and  $m$  respectively.  $A, x, y$  are complex arrays. When CVPRD1 is called,  $Ax + y$  is computed and stored in  $y$ .

**Remark.** It is assumed that the arrays  $A, x, y$  do not overlap.

**Programmer.** A. H. Morris

**CALL CTPRD( $m, n, A, IA, JA, x, y$ )**

The argument  $x$  is a row vector of dimension  $m$  and  $y$  an array of dimension  $n$ .  $A, x, y$  are complex arrays. When CTPRD is called,  $xA$  is computed and stored in  $y$ .

**Remark.** It is assumed that the arrays  $A, x, y$  do not overlap.

**Programmer.** A. H. Morris

**CALL CTPRD1( $m, n, A, IA, JA, x, y$ )**

The arguments  $x$  and  $y$  are row vectors of dimension  $m$  and  $n$  respectively.  $A, x, y$  are complex arrays. When CTPRD1 is called,  $xA + y$  is computed and stored in  $y$ .

**Remark.** It is assumed that the arrays  $A, x, y$  do not overlap.

**Programmer.** A. H. Morris



## ORDERING THE ROWS OF A SPARSE MATRIX BY INCREASING LENGTH

Let  $A$  be a sparse  $m \times n$  matrix stored in the arrays  $A, IA, JA$ . The following subroutine is available for ordering the rows of the matrix by increasing length.

**CALL SPORD**( $m, n, IA, R, IWK$ )

$R$  is an integer array of dimension  $m$ . When SPORD is called, the rows of the matrix are ordered by increasing length. The row ordering is given in  $R$ .

$IWK$  is an integer array of dimension  $m + n + 1$  or larger that is used for a work space.

**Remark.** If rows  $i_1, \dots, i_k$  are the rows of length  $\ell$ , then the indices  $i_1, \dots, i_k$  are listed in  $R$  in increasing sequence.

**Programmer.** A. H. Morris

## REORDERING SPARSE MATRICES INTO BLOCK TRIANGULAR FORM

Let  $A$  be a sparse  $n \times n$  matrix stored in the arrays  $A, IA, JA$ . Then the subroutine BLKORD is available for reordering the rows and columns of  $A$  so that one has a lower block triangular matrix

$$(*) \quad \begin{pmatrix} A_{11} & & & 0 \\ A_{21} & A_{22} & & \\ \vdots & \vdots & \ddots & \\ A_{k1} & A_{k2} & \dots & A_{kk} \end{pmatrix}$$

where the blocks  $A_{ii}$  are square and cannot themselves be reordered into lower block triangular form.

**CALL BLKORD**( $n, IA, JA, R, C, IB, k, IWK, IERR$ )

$R$  and  $C$  are integer arrays of dimension  $n$ , and  $IERR$  is an integer variable. When BLKORD is called, the rows of the matrix are first ordered so that the main diagonal contains a maximum number of nonzeros. After this is done then  $IERR$  = the number of zeros that appear on the diagonal. If  $IERR = 0$  then the rows and columns of the matrix are ordered into block triangular form (\*). The row ordering is given in  $R$  and the column ordering is given in  $C$ .

$IB$  is an integer array of dimension  $n$  and  $k$  is an integer variable. When the matrix has been ordered into block triangular form (\*) then  $k$  = the number of blocks  $A_{ii}$ . Also  $IB(i) =$  the row number in the block triangular matrix of the beginning of block  $A_{ii}$  ( $i = 1, \dots, k$ ).

$IWK$  is an integer array of dimension  $5n$  or larger that is used for a work space by the routine.

**Error Return.** If  $IERR \neq 0$  then the routine terminates. In this case,  $R$  contains the row ordering that gives the main diagonal with the maximum number of nonzeros.

**Remark.**  $IA, JA$ , and  $n$  are not modified by the routine.

**Programming.** BLKORD employs the subroutines MC21A, MC21B, MC13D, MC13E designed by I. S. Duff and J. K. Reid (AERE Harwell, England).

### References.

- (1) Duff, I. S., "On Algorithms for Obtaining a Maximum Transversal," *ACM Trans. Math Software* 7 (1981), pp.315-330.
- (2) Duff, I. S. and Reid, J. K., "An Implementation of Tarjan's Algorithm for the Block Triangularization of a Matrix," *ACM Trans. Math Software* 4 (1978), pp. 137-147.

## SOLUTION OF SPARSE SYSTEMS OF REAL LINEAR EQUATIONS

Let  $A$  be a nonsingular  $n \times n$  sparse real matrix stored in the arrays  $A, IA, JA$  and let  $b$  be a real column vector of dimension  $n$ . The subroutines SPSLV and RSLV are available for solving the system of equations  $Ax = b$ , and the subroutine TSLV is available for solving the transposed system of equations  $A^t x = b$ . These routines employ partial pivot gauss elimination with column interchanges to first obtain an LU decomposition of  $A$ . If SPSLV is called then only the off-diagonal nonzero elements of  $U$  are stored, and then the equations are solved. However, if RSLV or TSLV is called then the off-diagonal nonzero elements of both  $L$  and  $U$  are stored. Thus RSLV and TSLV will frequently require at least double the amount of storage needed by SPSLV, but they can be recalled to solve other systems of equations  $Ax = r$  and  $A^t x = r$  without having to redecompose the matrix  $A$ . Moreover, since RSLV and TSLV will always generate the same LU decomposition of  $A$ , RSLV can be called to decompose  $A$  and solve a system of equations  $Ax = b$ , and then TSLV can be called to solve a transposed system of equations  $A^t x = r$  using the decomposition obtained by RSLV.

**CALL SPSLV( $n, A, IA, JA, b, R, C, MAX, X, IWK, WK, IERR$ )**

It is assumed that  $n \geq 1$  and that  $X$  is an array of dimension  $n$ . The solution of the system of equations  $Ax = b$  is computed and stored in  $X$ .  $A, IA, JA$  and  $b$  are not modified by the routine.

$R$  is an integer array of  $n$  entries specifying the order in which the  $n$  rows of  $A$  are to be examined and processed. For example, if  $R$  contains the entries  $i_1, \dots, i_n$  then the algorithm first performs operations on row  $i_1$ , next on row  $i_2$ , etc. It is well known that the order in which the rows of a sparse matrix are processed can have a significant impact on the overall performance of a subroutine such as SPSLV. Thus  $R$  must be chosen judiciously.  $R$  is not modified by the routine.

$C$  is an integer array of  $n$  entries which plays a role similar to  $R$ . On input,  $C$  specifies a suggested order in which the  $n$  columns of  $A$  should be ordered for selection of the pivot elements. For example, if  $C$  contains the entries  $j_1, \dots, j_n$  then it is suggested that the first pivot element may be from column  $j_1$ , the second pivot element from column  $j_2$ , etc. However, since partial pivoting with column interchange is performed, on output  $C$  may have been modified. On output,  $C$  will contain the actual order of the  $n$  columns from which the pivot elements were selected. This order will depend on  $A$  and  $R$ , and not on  $b$ .

$IWK$  and  $WK$  are arrays for internal use by the routine, and  $MAX$  is an input argument. When SPSLV is called, an LU decomposition of  $A$  is first obtained where  $U$  is a unit upper triangular matrix. The off-diagonal portion of  $U$  is stored in sparse form in  $IWK$  and  $WK$ .  $MAX$  is an estimate of the maximum number of off-diagonal elements of  $U$  that might be nonzero and have to be stored ( $MAX \leq n(n-1)/2$ ).  $IWK$  is an integer array of dimension  $3n + MAX + 2$  or larger, and  $WK$  is a real array of dimension  $n + MAX$  or larger.

IERR is an integer variable that reports the status of the results. When the routine terminates, IERR has one of the following values:

- IERR > 0       $Ax = b$  was solved. IERR =  $\max\{1, m\}$  where  $m$  is the number of off-diagonal nonzero elements of  $U$ .
- IERR = 0      The argument  $n$  is nonpositive.
- IERR =  $-k$       Row  $R(k)$  of  $A$  is null.
- IERR =  $-n - k$       Row  $R(k)$  of  $A$  has a duplicate entry.
- IERR =  $-2n - k$       Row  $R(k)$  of  $A$  has been reduced to a row containing only zeros.
- IERR =  $-3n - k$       Row  $k$  of the upper triangular matrix exceeds storage. MAX must be increased.

When an error is detected, the routine immediately terminates.

**Remarks.** The amount of storage MAX depends critically on the row ordering given in  $R$ . If it is suspected that the rows and columns of  $A$  can be reordered so that one has a lower block triangular matrix

$$\begin{pmatrix} A_{11} & & & 0 \\ A_{21} & A_{22} & & \\ \vdots & \vdots & \ddots & \\ A_{k1} & A_{k2} & \dots & A_{kk} \end{pmatrix}$$

then the subroutine BLKORD should first be tried. This subroutine will specify an ordering for lower block triangular form if one exists. However, if such an ordering does not exist and one is uncertain how to order the rows, then the row ordering given by the subroutine SPORD frequently yields good results. In any case, the selection of an initial column ordering  $C$  is never bothersome since partial pivoting with column interchanges is performed. The initial ordering  $C(i) = i$  ( $i = 1, \dots, n$ ) always suffices.

**Programming.** SPSLV is a modification by A. H. Morris of the subroutine NSPIV. SPSLV employs the subroutine NSPIV1. NSPIV and NSPIV1 were written by Andrew H. Sherman (University of Texas at Austin).

**Reference.** Sherman, Andrew H., "Algorithms for Sparse Gaussian Elimination with Partial Pivoting," *ACM Trans. Math Software* 4 (1978), pp.330-338.

CALL RSLV(MO,n,A,IA,JA,b,R,C,MAX,X,IWK,WK,IERR)  
CALL TSLV(MO,n,A,IA,JA,b,R,C,MAX,X,IWK,WK,IERR)

RSLV is called for solving  $Ax = b$ , and TSLV is called for solving  $A^t x = b$ . MO is an input argument which specifies if RSLV or TSLV is being called for the first time. On an initial call, MO = 0 and we have the following setup:

It is assumed that  $n \geq 1$  and that  $X$  is an array of dimension  $n$ . The solution of the system of equations is stored in  $X$ .  $A$ ,  $IA$ ,  $JA$  are not modified by the routines.  $X$  and  $b$  may share the same storage area. If  $X$  is a separate storage area then  $b$  is not modified by the routines.

$R$  is an integer array of  $n$  entries specifying the order in which the  $n$  rows of  $A$  are to be examined and processed. For example, if  $R$  contains the entries  $i_1, \dots, i_n$  then the algorithm first performs operations on row  $i_1$ , next on row  $i_2$ , etc. It is well known that the order in which the rows of a sparse matrix are processed can have a significant impact on the overall performance of subroutines such as RSLV and TSLV. Thus  $R$  must be chosen judiciously.  $R$  is not modified by the routine.

$C$  is an integer array of  $n$  entries which plays a role similar to  $R$ . On input,  $C$  specifies a suggested order in which the  $n$  columns of  $A$  should be ordered for selection of the pivot elements. For example, if  $C$  contains the entries  $j_1, \dots, j_n$  then it is suggested that the first pivot element may be from column  $j_1$ , the second pivot element from column  $j_2$ , etc. However, since partial pivoting with column interchange is performed, on output  $C$  may have been modified. On output,  $C$  will contain the actual order of the  $n$  columns from which the pivot elements were selected. This order will depend on  $A$  and  $R$ , and not on  $b$ .

IWK and WK are arrays for internal use by the routines, and MAX is an input argument. On an initial call to RSLV or TSLV, an LU decomposition of  $A$  is first obtained where  $L$  is a lower triangular matrix and  $U$  a unit upper triangular matrix. The off-diagonal portions of  $L$  and  $U$  are stored in sparse form in IWK and WK. MAX is an estimate of the maximum number of off-diagonal elements of  $L$  and  $U$  that might be nonzero and have to be stored ( $\text{MAX} \leq n(n-1)$ ). IWK is an integer array of dimension  $4n + \text{MAX} + 2$  or larger, and WK is a real array of dimension  $2n + \text{MAX}$  or larger.

On an initial call to RSLV or TSLV, IERR is an integer variable that reports the status of the results. When the routine terminates, IERR has one of the following values:

- |                   |                                                                                                                                         |
|-------------------|-----------------------------------------------------------------------------------------------------------------------------------------|
| IERR > 0          | The system of equations was solved. IERR = max{1, $m$ } where $m$ is the total number of off-diagonal nonzero elements of $L$ and $U$ . |
| IERR = 0          | The argument $n$ is nonpositive.                                                                                                        |
| IERR = - $k$      | Row $R(k)$ of $A$ is null.                                                                                                              |
| IERR = - $n - k$  | Row $R(k)$ of $A$ has a duplicate entry.                                                                                                |
| IERR = - $2n - k$ | Row $R(k)$ of $A$ has been reduced to a row containing only zeros.                                                                      |
| IERR = - $3n - k$ | Row $k$ of $L$ or $U$ exceeds storage. MAX must be increased.                                                                           |

When an error is detected, the routine immediately terminates.

After an initial call to RSLV or TSLV, if IERR > 0 on output then either routine may be called with  $\text{MO} \neq 0$ . When  $\text{MO} \neq 0$  it is assumed that only  $b$  may have been modified. RSLV is called for solving the new set of equations  $Ax = b$ , and TSLV is called for solving the new set of equations  $A^t x = b$ . The routine employs the LU decomposition obtained on the initial call to RSLV or TSLV to solve the new system of equations. The solution is stored in  $X$ . As before,  $X$  and  $b$  may share the same storage area. If  $\text{MO} \neq 0$  then only  $n$ ,  $R$ ,  $C$ , IWK, and WK are used.  $A$ ,  $IA$ ,  $JA$ , MAX, and IERR are not referenced by the routine.

**Note.** The remarks concerning the ordering of the rows and columns of  $A$  when SPSLV is

used hold also for RSLV and TSLV.

**Programming.** RSLV calls the subroutines RSLV1 and SPLU, and TSLV calls the subroutines TSLV1 and SPLU. These routines were written by A. H. Morris.

## DOUBLE PRECISION SOLUTION OF SPARSE SYSTEMS OF REAL LINEAR EQUATIONS

Let  $A$  be a nonsingular  $n \times n$  sparse double precision matrix stored in the arrays  $A, IA, JA$  and let  $b$  be a double precision column vector of dimension  $n$ . The subroutines DSPSLV and DSLV are available for solving the system of equations  $Ax = b$ , and the subroutine DTSLV is available for solving the transposed system of equations  $A^t x = b$ . These routines employ partial pivot Gauss elimination with column interchanges to first obtain an LU decomposition of  $A$ . If DSPSLV is called then only the off-diagonal nonzero elements of  $U$  are stored, and then the equations are solved. However, if DSLV or DTSLV is called then the off-diagonal nonzero elements of both  $L$  and  $U$  are stored. Thus DSLV and DTSLV will frequently require at least double the amount of storage needed by DSPSLV, but they can be recalled to solve other systems of equations  $Ax = r$  and  $A^t x = r$  without having to redecompose the matrix  $A$ . Moreover, since DSLV and DTSLV will always generate the same LU decomposition of  $A$ , DSLV can be called to decompose  $A$  and solve a system of equations  $Ax = b$ , and then DTSLV can be called to solve a transposed system of equations  $A^t x = r$  using the decomposition obtained by DSLV.

**CALL DSPSLV( $n, A, IA, JA, b, R, C, MAX, X, IWK, WK, IERR$ )**

$A, b$ , and  $X$  are double precision arrays. It is assumed that  $n \geq 1$  and that  $X$  is an array of dimension  $n$ . The solution of the system of equations  $Ax = b$  is computed and stored in  $X$ .  $A, IA, JA$  and  $b$  are not modified by the routine.

$R$  is an integer array of  $n$  entries specifying the order in which the  $n$  rows of  $A$  are to be examined and processed. For example, if  $R$  contains the entries  $i_1, \dots, i_n$  then the algorithm first performs operations on row  $i_1$ , next on row  $i_2$ , etc. It is well known that the order in which the rows of a sparse matrix are processed can have a significant impact on the overall performance of a subroutine such as DSPSLV. Thus  $R$  must be chosen judiciously.  $R$  is not modified by the routine.

$C$  is an integer array of  $n$  entries which plays a role similar to  $R$ . On input,  $C$  specifies a suggested order in which the  $n$  columns of  $A$  should be ordered for selection of the pivot elements. For example, if  $C$  contains the entries  $j_1, \dots, j_n$  then it is suggested that the first pivot element may be from column  $j_1$ , the second pivot element from column  $j_2$ , etc. However, since partial pivoting with column interchange is performed, on output  $C$  may have been modified. On output,  $C$  will contain the actual order of the  $n$  columns from which the pivot elements were selected. This order will depend on  $A$  and  $R$ , and not on  $b$ .

$IWK$  and  $WK$  are arrays for internal use by the routine, and  $MAX$  is an input argument. When DSPSLV is called, an LU decomposition of  $A$  is first obtained where  $U$  is a unit upper triangular matrix. The off-diagonal portion of  $U$  is stored in sparse form in  $IWK$  and  $WK$ .  $MAX$  is an estimate of the maximum number of off-diagonal elements of  $U$  that might be nonzero and have to be stored ( $MAX \leq n(n-1)/2$ ).  $IWK$  is an integer array of dimension  $3n + MAX + 2$  or larger, and  $WK$  is a double precision array of dimension  $n + MAX$  or larger.

IERR is an integer variable that reports the status of the results. When the routine terminates, IERR has one of the following values:

- IERR > 0       $Ax = b$  was solved. IERR =  $\max\{1, m\}$  where  $m$  is the number of off-diagonal nonzero elements of  $U$ .
- IERR = 0      The argument  $n$  is nonpositive.
- IERR =  $-k$       Row  $R(k)$  of  $A$  is null.
- IERR =  $-n - k$       Row  $R(k)$  of  $A$  has a duplicate entry.
- IERR =  $-2n - k$       Row  $R(k)$  of  $A$  has been reduced to a row containing only zeros.
- IERR =  $-3n - k$       Row  $k$  of the upper triangular matrix exceeds storage. MAX must be increased.

When an error is detected, the routine immediately terminates.

**Remarks.** The amount of storage MAX depends critically on the row ordering given in  $R$ . If it is suspected that the rows and columns of  $A$  can be reordered so that one has a lower block triangular matrix

$$\begin{pmatrix} A_{11} & & & 0 \\ A_{21} & A_{22} & & \\ \vdots & \vdots & \ddots & \\ A_{k1} & A_{k2} & \dots & A_{kk} \end{pmatrix}$$

then the subroutine BLKORD should first be tried. This subroutine will specify an ordering for lower block triangular form if one exists. However, if such an ordering does not exist and one is uncertain how to order the rows, then the row ordering given by the subroutine SPORD frequently yields good results. In any case, the selection of an initial column ordering  $C$  is never bothersome since partial pivoting with column interchanges is performed. The initial ordering  $C(i) = i$  ( $i = 1, \dots, n$ ) always suffices.

**Programming.** DSPSLV an adaptation by A. H. Morris of the subroutine NSPIV written by Andrew H. Sherman (University of Texas at Austin). DSPSLV employs the subroutine DNSPIV.

**Reference.** Sherman, Andrew H., "Algorithms for Sparse Gaussian Elimination with Partial Pivoting," *ACM Trans. Math Software* 4 (1978), pp.330-338.

CALL DSLV(MO,n,A,IA,JA,b,R,C,MAX,X,IWK,WK,IERR)  
CALL DTSLV(MO,n,A,IA,JA,b,R,C,MAX,X,IWK,WK,IERR)

DSLV is called for solving  $Ax = b$ , and DTSLV is called for solving  $A^t x = b$ . MO is an input argument which specifies if DSLV or DTSLV is being called for the first time. On an initial call, MO = 0 and we have the following setup:

$A$ ,  $b$ , and  $X$  are double precision arrays. It is assumed that  $n \geq 1$  and that  $X$  is an array of dimension  $n$ . The solution of the system of equations is stored in  $X$ .  $A$ ,  $IA$ ,  $JA$  are not modified by the routines.  $X$  and  $b$  may share the same storage area. If  $X$  is a separate storage area then  $b$  is not modified by the routines.



$R$  is an integer array of  $n$  entries specifying the order in which the  $n$  rows of  $A$  are to be examined and processed. For example, if  $R$  contains the entries  $i_1, \dots, i_n$  then the algorithm first performs operations on row  $i_1$ , next on row  $i_2$ , etc. It is well known that the order in which the rows of a sparse matrix are processed can have a significant impact on the overall performance of subroutines such as DSLV and DTSLV. Thus  $R$  must be chosen judiciously.  $R$  is not modified by the routine.

$C$  is an integer array of  $n$  entries which plays a role similar to  $R$ . On input,  $C$  specifies a suggested order in which the  $n$  columns of  $A$  should be ordered for selection of the pivot elements. For example, if  $C$  contains the entries  $j_1, \dots, j_n$  then it is suggested that the first pivot element may be from column  $j_1$ , the second pivot element from column  $j_2$ , etc. However, since partial pivoting with column interchange is performed, on output  $C$  may have been modified. On output,  $C$  will contain the actual order of the  $n$  columns from which the pivot elements were selected. This order will depend on  $A$  and  $R$ , and not on  $b$ .

IWK and WK are arrays for internal use by the routines, and MAX is an input argument. On an initial call to DSLV or DTSLV, an LU decomposition of  $A$  is first obtained where  $L$  is a lower triangular matrix and  $U$  a unit upper triangular matrix. The off-diagonal portions of  $L$  and  $U$  are stored in sparse form in IWK and WK. MAX is an estimate of the maximum number of off-diagonal elements of  $L$  and  $U$  that might be nonzero and have to be stored ( $\text{MAX} \leq n(n-1)$ ). IWK is an integer array of dimension  $4n + \text{MAX} + 2$  or larger, and WK is a double precision array of dimension  $2n + \text{MAX}$  or larger.

On an initial call to DSLV or DTSLV, IERR is an integer variable that reports the status of the results. When the routine terminates, IERR has one of the following values:

- |                  |                                                                                                                                         |
|------------------|-----------------------------------------------------------------------------------------------------------------------------------------|
| IERR > 0         | The system of equations was solved. IERR = max{1, $m$ } where $m$ is the total number of off-diagonal nonzero elements of $L$ and $U$ . |
| IERR = 0         | The argument $n$ is nonpositive.                                                                                                        |
| IERR = $-k$      | Row $R(k)$ of $A$ is null.                                                                                                              |
| IERR = $-n - k$  | Row $R(k)$ of $A$ has a duplicate entry.                                                                                                |
| IERR = $-2n - k$ | Row $R(k)$ of $A$ has been reduced to a row containing only zeros.                                                                      |
| IERR = $-3n - k$ | Row $k$ of $L$ or $U$ exceeds storage. MAX must be increased.                                                                           |

When an error is detected, the routine immediately terminates.

After an initial call to DSLV or DTSLV, if IERR > 0 on output then either routine may be called with  $\text{MO} \neq 0$ . When  $\text{MO} \neq 0$  it is assumed that only  $b$  may have been modified. DSLV is called for solving the new set of equations  $Ax = b$ , and DTSLV is called for solving the new set of equations  $A^t x = b$ . The routine employs the LU decomposition obtained on the initial call to DSLV or DTSLV to solve the new system of equations. The solution is stored in  $X$ . As before,  $X$  and  $b$  may share the same storage area. If  $\text{MO} \neq 0$  then only  $n$ ,  $R$ ,  $C$ , IWK, and WK are used.  $A$ ,  $IA$ ,  $JA$ , MAX, and IERR are not referenced by the routine.

**Note.** The remarks concerning the ordering of the rows and columns of  $A$  when DSPSLV

is used hold also for DSLV and DTSLV.

**Programming.** DSLV calls the subroutines DSLV1 and DSPLU, and DTSLV calls the subroutines DTSLV1 and DSPLU. These routines were written by A. H. Morris.

## SOLUTION OF SPARSE SYSTEMS OF COMPLEX LINEAR EQUATIONS

Let  $A$  be a nonsingular  $n \times n$  sparse complex matrix stored in the arrays  $A, IA, JA$  and let  $b$  be a complex column vector of dimension  $n$ . The subroutines CSPSLV and CSLV are available for solving the system of equations  $Ax = b$ , and the subroutine CTSLV is available for solving the transposed system of equations  $A^t x = b$ . These routines employ partial pivot Gauss elimination with column interchanges to first obtain an LU decomposition of  $A$ . If CSPSLV is called then only the off-diagonal nonzero elements of  $U$  are stored, and then the equations are solved. However, if CSLV or CTSLV is called then the off-diagonal nonzero elements of both  $L$  and  $U$  are stored. Thus CSLV and CTSLV will frequently require at least double the amount of storage needed by CSPSLV, but they can be recalled to solve other systems of equations  $Ax = r$  and  $A^t x = r$  without having to redecompose the matrix  $A$ . Moreover, since CSLV and CTSLV will always generate the same LU decomposition of  $A$ , CSLV can be called to decompose  $A$  and solve a system of equations  $Ax = b$ , and then CTSLV can be called to solve a transposed system of equations  $A^t x = r$  using the decomposition obtained by CSLV.

**CALL CSPSLV( $n, A, IA, JA, b, R, C, MAX, X, IWK, WK, IERR$ )**

$A, b$ , and  $X$  are complex arrays. It is assumed that  $n \geq 1$  and that  $X$  is an array of dimension  $n$ . The solution of the system of equations  $Ax = b$  is computed and stored in  $X$ .  $A, IA, JA$  and  $b$  are not modified by the routine.

$R$  is an integer array of  $n$  entries specifying the order in which the  $n$  rows of  $A$  are to be examined and processed. For example, if  $R$  contains the entries  $i_1, \dots, i_n$  then the algorithm first performs operations on row  $i_1$ , next on row  $i_2$ , etc. It is well known that the order in which the rows of a sparse matrix are processed can have a significant impact on the overall performance of a subroutine such as CSPSLV. Thus  $R$  must be chosen judiciously.  $R$  is not modified by the routine.

$C$  is an integer array of  $n$  entries which plays a role similar to  $R$ . On input,  $C$  specifies a suggested order in which the  $n$  columns of  $A$  should be ordered for selection of the pivot elements. For example, if  $C$  contains the entries  $j_1, \dots, j_n$  then it is suggested that the first pivot element may be from column  $j_1$ , the second pivot element from column  $j_2$ , etc. However, since partial pivoting with column interchange is performed, on output  $C$  may have been modified. On output,  $C$  will contain the actual order of the  $n$  columns from which the pivot elements were selected. This order will depend on  $A$  and  $R$ , and not on  $b$ .

$IWK$  and  $WK$  are arrays for internal use by the routine, and  $MAX$  is an input argument. When CSPSLV is called, an LU decomposition of  $A$  is first obtained where  $U$  is a unit upper triangular matrix. The off-diagonal portion of  $U$  is stored in sparse form in  $IWK$  and  $WK$ .  $MAX$  is an estimate of the maximum number of off-diagonal elements of  $U$  that might be nonzero and have to be stored ( $MAX \leq n(n-1)/2$ ).  $IWK$  is an integer array of dimension  $3n + MAX + 2$  or larger, and  $WK$  is a complex array of dimension  $n + MAX$  or larger.

IERR is an integer variable that reports the status of the results. When the routine terminates, IERR has one of the following values:

- IERR > 0       $Ax = b$  was solved. IERR =  $\max\{1, m\}$  where  $m$  is the number of off-diagonal nonzero elements of  $U$ .
- IERR = 0      The argument  $n$  is nonpositive.
- IERR =  $-k$       Row  $R(k)$  of  $A$  is null.
- IERR =  $-n - k$       Row  $R(k)$  of  $A$  has a duplicate entry.
- IERR =  $-2n - k$       Row  $R(k)$  of  $A$  has been reduced to a row containing only zeros.
- IERR =  $-3n - k$       Row  $k$  of the upper triangular matrix exceeds storage. MAX must be increased.

When an error is detected, the routine immediately terminates.

**Remarks.** The amount of storage MAX depends critically on the row ordering given in  $R$ . If it is suspected that the rows and columns of  $A$  can be reordered so that one has a lower block triangular matrix

$$\begin{pmatrix} A_{11} & & & 0 \\ A_{21} & A_{22} & & \\ \vdots & \vdots & \ddots & \\ A_{k1} & A_{k2} & \dots & A_{kk} \end{pmatrix}$$

then the subroutine BLKORD should first be tried. This subroutine will specify an ordering for lower block triangular form if one exists. However, if such an ordering does not exist and one is uncertain how to order the rows, then the row ordering given by the subroutine SPORD frequently yields good results. In any case, the selection of an initial column ordering  $C$  is never bothersome since partial pivoting with column interchanges is performed. The initial ordering  $C(i) = i$  ( $i = 1, \dots, n$ ) always suffices.

**Programming.** CSPSLV is an adaptation by A. H. Morris of the subroutine NSPIV written by Andrew H. Sherman (University of Texas at Austin). CSPSLV employs the subroutine CNSPIV.

**Reference.** Sherman, Andrew H., "Algorithms for Sparse Gaussian Elimination with Partial Pivoting," *ACM Trans. Math Software* 4 (1978), pp.330-338.

CALL CSLV(MO,n,A,IA,JA,b,R,C,MAX,X,IWK,WK,IERR)  
CALL CTSLV(MO,n,A,IA,JA,b,R,C,MAX,X,IWK,WK,IERR)

CSLV is called for solving  $Ax = b$ , and CTSLV is called for solving  $A^t x = b$ . MO is an input argument which specifies if CSLV or CTSLV is being called for the first time. On an initial call, MO = 0 and we have the following setup:

$A$ ,  $b$ , and  $X$  are complex arrays. It is assumed that  $n \geq 1$  and that  $X$  is an array of dimension  $n$ . The solution of the system of equations is stored in  $X$ .  $A$ ,  $IA$ ,  $JA$  are not modified by the routines.  $X$  and  $b$  may share the same storage area. If  $X$  is a separate storage area then  $b$  is not modified by the routines.

$R$  is an integer array of  $n$  entries specifying the order in which the  $n$  rows of  $A$  are to be examined and processed. For example, if  $R$  contains the entries  $i_1, \dots, i_n$  then the algorithm first performs operations on row  $i_1$ , next on row  $i_2$ , etc. It is well known that the order in which the rows of a sparse matrix are processed can have a significant impact on the overall performance of subroutines such as CSLV and CTSVL. Thus  $R$  must be chosen judiciously.  $R$  is not modified by the routine.

$C$  is an integer array of  $n$  entries which plays a role similar to  $R$ . On input,  $C$  specifies a suggested order in which the  $n$  columns of  $A$  should be ordered for selection of the pivot elements. For example, if  $C$  contains the entries  $j_1, \dots, j_n$  then it is suggested that the first pivot element may be from column  $j_1$ , the second pivot element from column  $j_2$ , etc. However, since partial pivoting with column interchange is performed, on output  $C$  may have been modified. On output,  $C$  will contain the actual order of the  $n$  columns from which the pivot elements were selected. This order will depend on  $A$  and  $R$ , and not on  $b$ .

IWK and WK are arrays for internal use by the routines, and MAX is an input argument. On an initial call to CSLV or CTSVL, an LU decomposition of  $A$  is first obtained where  $L$  is a lower triangular matrix and  $U$  a unit upper triangular matrix. The off-diagonal portions of  $L$  and  $U$  are stored in sparse form in IWK and WK. MAX is an estimate of the maximum number of off-diagonal elements of  $L$  and  $U$  that might be nonzero and have to be stored ( $\text{MAX} \leq n(n-1)$ ). IWK is an integer array of dimension  $4n + \text{MAX} + 2$  or larger, and WK is a complex array of dimension  $2n + \text{MAX}$  or larger.

On an initial call to CSLV or CTSVL, IERR is an integer variable that reports the status of the results. When the routine terminates, IERR has one of the following values:

- |                   |                                                                                                                                         |
|-------------------|-----------------------------------------------------------------------------------------------------------------------------------------|
| IERR > 0          | The system of equations was solved. IERR = max{1, $m$ } where $m$ is the total number of off-diagonal nonzero elements of $L$ and $U$ . |
| IERR = 0          | The argument $n$ is nonpositive.                                                                                                        |
| IERR = - $k$      | Row $R(k)$ of $A$ is null.                                                                                                              |
| IERR = - $n - k$  | Row $R(k)$ of $A$ has a duplicate entry.                                                                                                |
| IERR = - $2n - k$ | Row $R(k)$ of $A$ has been reduced to a row containing only zeros.                                                                      |
| IERR = - $3n - k$ | Row $k$ of $L$ or $U$ exceeds storage. MAX must be increased.                                                                           |

When an error is detected, the routine immediately terminates.

After an initial call to CSLV or CTSVL, if IERR > 0 on output then either routine may be called with  $\text{MO} \neq 0$ . When  $\text{MO} \neq 0$  it is assumed that only  $b$  may have been modified. CSLV is called for solving the new set of equations  $Ax = b$ , and CTSVL is called for solving the new set of equations  $A^t x = b$ . The routine employs the LU decomposition obtained on the initial call to CSLV or CTSVL to solve the new system of equations. The solution is stored in  $X$ . As before,  $X$  and  $b$  may share the same storage area. If  $\text{MO} \neq 0$  then only  $n$ ,  $R$ ,  $C$ , IWK, and WK are used.  $A$ ,  $IA$ ,  $JA$ , MAX, and IERR are not referenced by the routine.

**Note.** The remarks concerning the ordering of the rows and columns of  $A$  when CPSLV is used hold also for CSLV and CTSLV.

**Programming.** CSLV calls the subroutines CSLV1 and CSPLU, and CTSLV calls the subroutines CTSLV1 and CSPLU. These routines were written by A. H. Morris.

## COMPUTATION OF EIGENVALUES OF GENERAL REAL MATRICES

The subroutines EIG and EIG1 are available for computing the eigenvalues of real matrices. These routines frequently yield results accurate to 13–14 significant digits. Indeed, for symmetric matrices they may give 2 or more digits better accuracy than the routines designed specifically for symmetric matrices. However, if the eigenvalues are not distinct or if they are exceedingly tightly clustered, then a severe drop in accuracy can occur when the matrix is not symmetric. In this case one should not expect more than 7–8 digit accuracy.

```
CALL EIG(IBAL,A,ka,n,WR,WI,IERR)
CALL EIG1(IBAL,A,ka,n,WR,WI,IERR)
```

$A$  is a matrix of order  $n \geq 1$  and  $WR, WI$  are real arrays of dimension  $n$  or larger. When EIG or EIG1 is called then the eigenvalues  $\lambda_1, \dots, \lambda_n$  of  $A$  are computed. The real parts of the eigenvalues are stored in  $WR(1), \dots, WR(n)$  and the imaginary parts are stored in  $WI(1), \dots, WI(n)$ . The eigenvalues are unordered except that complex conjugate pairs of eigenvalues appear consecutively with the eigenvalue having the positive imaginary part being first.

IBAL and  $ka$  are input arguments. The argument  $ka$  is the number of rows in the dimension statement for  $A$  in the calling program. IBAL may be any integer. If  $IBAL \neq 0$  then the routines balance  $A$  before they compute the eigenvalues. Otherwise, if  $IBAL = 0$  then  $A$  is not balanced.

**Error Return.** IERR is an integer variable. If all the eigenvalues are found then IERR is set to 0. Otherwise, if more than 30 iterations are required to compute the  $j^{th}$  eigenvalue  $\lambda_j$ , then IERR is set to  $j$  and the routine terminates. In this case, if  $j < n$  then the eigenvalues  $\lambda_{j+1}, \dots, \lambda_n$  will have been computed and the results stored in the  $WR$  and  $WI$  arrays.

### Remarks.

- (1) Even though the balancing operation does not increase the theoretical bounds on the errors, nevertheless at times it may result in a slight loss of accuracy. On the other hand, balancing requires little additional time and in certain cases can improve the accuracy by as much as 5–6 significant digits. Thus it is recommended that balancing be done.
- (2)  $A$  is destroyed during computation. EIG and EIG1 reduce  $A$  to upper Hessenberg form and then apply the QR algorithm to obtain the eigenvalues. They differ only in the choice of transformations used to reduce  $A$  to upper Hessenberg form. EIG employs elementary similarity transformations and EIG1 employs orthogonal similarity transformations. In theory the use of orthogonal transformations assures one of a tighter bound on the errors. However, since in practice matrices infrequently arise for which the orthogonal transformations actually generate more accurate results, and since the orthogonal transformations normally require more time than the elementary transformations, therefore EIG is the recommended routine.

**Programming.** EIG and EIG1 are driver routines for the EISPACK subroutines BALANC, ELMHS0, ORTHES, and HQR. These subroutines were developed at Argonne National Laboratory. The functions SPMPAR and IPMPAR are also used.

**Reference.** Smith, B. T., Boyle, J. M., et al., *Matrix Eigensystem Routines - EISPACK Guide* (Second Edition), Springer-Verlag, 1976.



## COMPUTATION OF EIGENVALUES AND EIGENVECTORS OF GENERAL REAL MATRICES

The subroutines EIGV and EIGV1 are available for computing the eigenvalues and eigenvectors of real matrices. These routines are extensions of the respective eigenvalue routines EIG and EIG1. Thus all comments made concerning the accuracy of the eigenvalues produced by EIG and EIG1 apply also to EIGV and EIGV1. In particular, EIGV and EIGV1 can frequently yield high precision results for the eigenvalues if they are distinct. However, be aware that errors in the eigenvalues, no matter how seemingly insignificant, can be considerably magnified in the computation of the eigenvectors. It is not at all unusual to obtain an eigenvalue and eigenvector where the eigenvalue is correct to within 2-3 units of the 14<sup>th</sup> significant digit, but the components of the corresponding eigenvector are only accurate to 9-10 significant digits. In the case of repeated eigenvalues the situation regarding the eigenvectors is totally unpredictable. The components of such an eigenvector may be correct to 6-7 significant digits, or the eigenvector may not even be an eigenvector! In this case the results should be checked.

```
CALL EIGV(IBAL, A, ka, n, WR, WI, ZR, ZI, IERR)
CALL EIGV1(IBAL, A, ka, n, WR, WI, ZR, ZI, IERR)
```

$A$  is a matrix of order  $n \geq 1$  and WR, WI are real arrays of dimension  $n$  or larger. When EIGV or EIGV1 is called the eigenvalues  $\lambda_1, \dots, \lambda_n$  and corresponding eigenvectors  $z_1, \dots, z_n$  are computed. The real parts of the eigenvalues are stored in WR(1), ..., WR( $n$ ) and the imaginary parts are stored in WI(1), ..., WI( $n$ ). The eigenvalues are unordered except that complex conjugate pairs of eigenvalues appear consecutively with the eigenvalue having the positive imaginary part being first.

The input argument  $ka$  is the number of rows in the dimension statement for  $A$  in the calling program. ZR and ZI are real arrays of dimension  $ka \times n$ . For  $j = 1, \dots, n$  the real parts of the components of the eigenvector  $z_j$  are stored in the  $j^{\text{th}}$  column of ZR (in locations ZR(1, $j$ ), ..., ZR( $n$ , $j$ )) and the imaginary parts are stored in the  $j^{\text{th}}$  column of ZI. The eigenvectors  $z_1, \dots, z_n$  are not normalized.

IBAL is an input argument that can be assigned any integer value. If IBAL  $\neq 0$  then the routines balance  $A$  before they compute the eigenvalues and eigenvectors. Otherwise, if IBAL = 0 then  $A$  is not balanced.

**Error Return.** IERR is an integer variable. If all the eigenvalues and eigenvectors are found then IERR is set to 0. Otherwise, if more than 30 iterations are required to compute the  $j^{\text{th}}$  eigenvalue  $\lambda_j$ , then IERR is set to  $j$  and the routine terminates. In this case, if  $j < n$  then the eigenvalues  $\lambda_{j+1}, \dots, \lambda_n$  will have been computed and the results stored in the WR and WI arrays. However, none of the eigenvectors will have been computed. The eigenvectors are computed only after all the eigenvalues have been obtained.

**Remarks.**

- (1) Even though the balancing operation does not increase the theoretical bounds on the errors, nevertheless at times it may result in a slight loss of accuracy. On the other hand, balancing requires little additional time and in certain cases can improve the accuracy of the eigenvalues by as much as 5-6 significant digits. When this occurs balancing will normally be needed to obtain the eigenvectors. In general, it is recommended that balancing be done.
- (2)  $A$  is destroyed during computation. EIGV and EIGV1 both reduce  $A$  to upper Hessenberg form, apply the QR algorithm to the Hessenberg matrix to obtain the eigenvalues, and then backsubstitute to generate the eigenvectors. They differ only in the choice of transformations used to reduce  $A$  to upper Hessenberg form. EIGV employs elementary similarity transformations and EIGV1 employs orthogonal similarity transformations. In theory the use of orthogonal transformations assures one of a tighter bound on the errors. However, since in practice matrices infrequently arise for which the orthogonal transformations actually generate more accurate results, and since the orthogonal transformations normally require more time than the elementary transformations, therefore EIGV is the recommended routine.

**Programming.** EIGV and EIGV1 are driver routines for the EISPACK subroutines BALANC, ELMHS0, ORTHES, ELTRN0, ORTRAN, HQR2, and BALBAK. These subroutines were developed at Argonne National Laboratory. The functions SPMPAR and IPMPAR are also used.

**Reference.** Smith, B. T., Boyle, J. M., et al., *Matrix Eigensystem Routines - EISPACK Guide* (Second Edition), Springer-Verlag, 1976.

## DOUBLE PRECISION COMPUTATION OF EIGENVALUES OF REAL MATRICES

The subroutine DEIG is available for the double precision computation of the eigenvalues of real matrices. This routine frequently yields results accurate to 26–28 significant digits. However, if the eigenvalues are not distinct or if they are exceedingly tightly clustered, then a severe drop in accuracy can occur. In this case one should not expect more than 13–14 digit accuracy.

**CALL DEIG**(IBAL, A, ka, n, WR, WI, IERR)

A is a double precision matrix of order  $n \geq 1$  and WR, WI are double precision arrays of dimension  $n$  or larger. When DEIG is called then the eigenvalues  $\lambda_1, \dots, \lambda_n$  of A are computed. The real parts of the eigenvalues are stored in WR(1), ..., WR(n) and the imaginary parts are stored in WI(1), ..., WI(n). The eigenvalues are unordered except that complex conjugate pairs of eigenvalues appear consecutively with the eigenvalue having the positive imaginary part being first.

IBAL and ka are input arguments. The argument ka is the number of rows in the dimension statement for A in the calling program. IBAL may be any integer. If IBAL  $\neq$  0 then the routine balances A before it computes the eigenvalues. Otherwise, if IBAL = 0 then A is not balanced.

**Error Return.** IERR is an integer variable. If all the eigenvalues are found then IERR is set to 0. Otherwise, if more than 50 iterations are required to compute the  $j^{\text{th}}$  eigenvalue  $\lambda_j$ , then IERR is set to  $j$  and the routine terminates. In this case, if  $j < n$  then the eigenvalues  $\lambda_{j+1}, \dots, \lambda_n$  will have been computed and the results stored in the WR and WI arrays.

### Remarks.

- (1) A is destroyed during computation.
- (2) DEIG is a double precision version of the eigenvalue routine EIG1.

**Programming.** DEIG is a driver routine for the subroutines DBAL, DORTH, and DHQR. These subroutines are double precision versions of the EISPACK subroutines BALANC, ORTHES, and HQR, developed at Argonne National Laboratory. The double precision versions were prepared by A. H. Morris. The functions DPMPAR and IPMPAR are also used.

**Reference.** Smith, B. T., Boyle, J. M., et al., *Matrix Eigensystem Routines – EISPACK Guide* (Second Edition), Springer-Verlag, 1976.

## DOUBLE PRECISION COMPUTATION OF EIGENVALUES AND EIGENVECTORS OF REAL MATRICES

The subroutine DEIGV is available for the double precision computation of the eigenvalues and eigenvectors of real matrices. This routine frequently yields values for the eigenvalues that are accurate to 26–28 significant digits. However, be aware that errors in the eigenvalues, no matter how seemingly insignificant, can be considerably magnified in the computation of the eigenvectors. If the eigenvalues are not distinct or if they are exceedingly tightly clustered, then a severe drop in accuracy can occur. In this case one should not expect the eigenvalues to have more than 13–14 digit accuracy.

**CALL DEIGV**(IBAL, A, ka, n, WR, WI, ZR, ZI, IERR)

A is a double precision matrix of order  $n \geq 1$  and WR, WI are double precision arrays of dimension  $n$  or larger. When DEIGV is called then the eigenvalues  $\lambda_1, \dots, \lambda_n$  and corresponding eigenvectors  $z_1, \dots, z_n$  are computed. The real parts of the eigenvalues are stored in WR(1), ..., WR( $n$ ) and the imaginary parts are stored in WI(1), ..., WI( $n$ ). The eigenvalues are unordered except that complex conjugate pairs of eigenvalues appear consecutively with the eigenvalue having the positive imaginary part being first.

The input argument ka is the number of rows in the dimension statement for A in the calling program. ZR and ZI are double precision arrays of dimension  $ka \times n$ . For  $j = 1, \dots, n$  the real parts of the components of the eigenvector  $z_j$  are stored in the  $j^{\text{th}}$  column of ZR (in locations ZR(1,  $j$ ), ..., ZR( $n$ ,  $j$ )) and the imaginary parts are stored in the  $j^{\text{th}}$  column of ZI. The eigenvectors  $z_1, \dots, z_n$  are not normalized.

IBAL is an input argument that can be assigned any integer value. If IBAL  $\neq 0$  then the routine balances A before it computes the eigenvalues and eigenvectors. Otherwise, if IBAL = 0 then A is not balanced.

**Error Return.** IERR is an integer variable. If all the eigenvalues and eigenvectors are found then IERR is set to 0. Otherwise, if more than 50 iterations are required to compute the  $j^{\text{th}}$  eigenvalue  $\lambda_j$ , then IERR is set to  $j$  and the routine terminates. In this case, if  $j < n$  then the eigenvalues  $\lambda_{j+1}, \dots, \lambda_n$  will have been computed and the results stored in the WR and WI arrays. However, none of the eigenvectors will have been computed. The eigenvectors are computed only after all the eigenvalues have been obtained.

### Remarks.

- (1) A is destroyed during computation.
- (2) DEIGV is a double precision version of the eigenvalue/eigenvector routine EIGV1.

**Programming.** DEIGV is a driver routine for the subroutines DBAL, DORTH, DORTRN, DHQR2, and DBABK. These subroutines are double precision versions of the EISPACK routines BALANC, ORTHES, ORTRAN, HQR2, and BALBAK, developed at Argonne National Laboratory. The double precision versions were prepared by A. H. Morris. The functions DPMPAR and IPMPAR are also used.

**Reference.** Smith, B. T., Boyle, J. M., et al., *Matrix Eigensystem Routines – EISPACK Guide* (Second Edition), Springer-Verlag, 1976.

## COMPUTATION OF EIGENVALUES OF SYMMETRIC REAL MATRICES

The subroutines SEIG and SEIG1 are available for computing the eigenvalues of symmetric real matrices. These routines frequently yield high precision results. SEIG is faster than SEIG1, but at times SEIG1 will produce better results when the symmetric matrix is tridiagonal. For arbitrary symmetric matrices it is not clear if there is any difference in the reliability of the routines.

```
CALL SEIG(A, ka, n, W, T, IERR)
CALL SEIG1(A, ka, n, W, T, IERR)
```

$A$  is a symmetric matrix of order  $n \geq 1$  and  $W$  an array of dimension  $n$  or larger. When SEIG or SEIG1 is called the eigenvalues  $\lambda_1, \dots, \lambda_n$  are computed and stored in  $W(1), \dots, W(n)$ . The eigenvalues are ordered so that  $\lambda_1 \leq \dots \leq \lambda_n$ .

$A$  may be packed or in standard form.<sup>1</sup> The input argument  $ka$  is a nonnegative integer. If  $ka = 0$  then  $A$  is assumed to be packed. Otherwise, if  $ka \neq 0$  then  $A$  is assumed to be in the standard format. In this case  $ka$  has the value:

$ka$  = the number of rows in the dimension statement for  $A$  in the calling program. It is assumed that  $ka \geq n$ . However, it is not required that  $A(i, j)$  be defined for  $i < j$ . Only the lower triangular elements of  $A$  are used.

$T$  is an array used for temporary storage. If SEIG is called then  $T$  must be of dimension  $2n$ . However, if SEIG1 is called then  $T$  need only be of dimension  $n$ .

**Error Return.** IERR is an integer variable. If all the eigenvalues are found then IERR is set to 0. Otherwise, if more than 30 iterations of the QL algorithm are required to compute the  $j^{\text{th}}$  eigenvalue  $\lambda_j$ , then IERR is set to  $j$ . In this case, if  $j > 1$  then the eigenvalues  $\lambda_1, \dots, \lambda_{j-1}$  will have been computed and stored in  $W$ . The eigenvalues are ordered so that  $\lambda_1 \leq \dots \leq \lambda_{j-1}$ . However, they need not be the smallest eigenvalues of  $A$ .

**Note.**  $A$  is destroyed during computation.

**Programming.** SEIG and SEIG1 are driver routines for the EISPACK subroutines TRED1, TRED3, TQLRAT, and IMTQL1. These subroutines were developed at Argonne National Laboratory. The function SPMPAR is also used.

**Reference.** Smith, B. T., Boyle, J. M., et al., *Matrix Eigensystem Routines - EISPACK Guide* (Second Edition), Springer-Verlag, 1976.

<sup>1</sup> For details on the packed format see the section on packing and unpacking symmetric matrices.

## COMPUTATION OF EIGENVALUES AND EIGENVECTORS OF SYMMETRIC REAL MATRICES

The subroutines SEIGV and SEIGV1 are available for computing the eigenvalues and eigenvectors of symmetric real matrices. These routines frequently yield high precision results for the eigenvalues. However, be aware that errors in the eigenvalues, no matter how seemingly insignificant, can be considerably magnified in the computation of the eigenvectors. It is not at all unusual to obtain an eigenvalue and eigenvector where the eigenvalue is correct to within 2-3 units of the 14<sup>th</sup> significant digit, but the components of the corresponding eigenvector are only accurate to 9-10 significant digits. SEIGV is faster than SEIGV1, but at times SEIGV1 will produce better results when the symmetric matrix is tridiagonal. For arbitrary symmetric matrices it is not clear if there is any difference in the reliability of the routines.

```
CALL SEIGV(A, ka, n, W, Z, T, IERR)
CALL SEIGV1(A, ka, n, W, Z, T, IERR)
```

$A$  is a symmetric matrix of order  $n \geq 1$  and  $W$  is an array of dimension  $n$  or larger. When SEIG or SEIGV1 is called the eigenvalues  $\lambda_1, \dots, \lambda_n$  and corresponding orthonormal eigenvectors  $z_1, \dots, z_n$  are computed. The eigenvalues are stored in  $W(1), \dots, W(n)$  and are ordered so that  $\lambda_1 \leq \dots \leq \lambda_n$ .

$A$  must be in the standard format, having the dimension  $ka \times n$ . It is assumed that  $ka \geq n$ . However, it is not required that  $A(i, j)$  be defined for  $i < j$ . Only the lower triangular elements of  $A$  are used.

$Z$  is an array of dimension  $ka \times n$  or larger. For  $j = 1, \dots, n$  the components of the eigenvector  $z_j$  are stored in the  $j^{\text{th}}$  column of  $Z$  (in locations  $Z(1, j), \dots, Z(n, j)$ ). To conserve memory one can let  $A$  and  $Z$  denote the same array.

$T$  is an array of dimension  $n$  used for temporary storage.

**Error Return.** IERR is an integer variable. If all the eigenvalues and eigenvectors are found then IERR is set to 0. Otherwise, if more than 30 iterations of the QL algorithm are required to compute the  $j^{\text{th}}$  eigenvalue  $\lambda_j$ , then IERR is set to  $j$ . In this case, if  $j > 1$  then the eigenvalues  $\lambda_1, \dots, \lambda_{j-1}$  and eigenvectors  $z_1, \dots, z_{j-1}$  will have been computed and stored in the  $W$  and  $Z$  arrays. However, the eigenvalues will be unordered.

**Note.**  $A$  is destroyed during computation.

**Programming.** SEIGV and SEIGV1 are driver routines for the EISPACK subroutines TRED2, TQL2, and IMTQL2. These subroutines were developed at Argonne National Laboratory. The function SPMPAR is also used.

**Reference.** Smith, B. T., Boyle, J. M., et al., *Matrix Eigensystem Routines - EISPACK Guide* (Second Edition), Springer-Verlag, 1976.

## COMPUTATION OF EIGENVALUES OF COMPLEX MATRICES

The subroutine CEIG is available for computing the eigenvalues of complex matrices. This routine frequently yields results accurate to 13–14 significant digits. However, if the eigenvalues are not distinct or if they are exceedingly tightly clustered, then a severe drop in accuracy can occur. In this case one should not expect more than 7–8 digit accuracy.

**CALL CEIG**(IBAL,AR,AI,ka,n,WR,WI,IERR)

AR and AI are real matrices of order  $n \geq 1$ , and WR and WI are real arrays of dimension  $n$  or larger. AR and AI are the real and imaginary portions of the complex matrix whose eigenvalues are to be computed. When CEIG is called the eigenvalues  $\lambda_1, \dots, \lambda_n$  are computed. The real parts of the eigenvalues are stored in WR(1), ..., WR( $n$ ) and the imaginary parts are stored in WI(1), ..., WI( $n$ ). The eigenvalues are unordered.

IBAL and ka are input arguments. It is assumed that ka is the number of rows in the dimension statements for AR and AI in the calling program. IBAL may be any integer. If IBAL  $\neq 0$  then the complex matrix (represented by AR and AI) is balanced before the eigenvalues are computed. Otherwise, if IBAL = 0 then the complex matrix is not balanced.

**Error Return.** IERR is an integer variable. If all the eigenvalues are found then IERR is set to 0. Otherwise, if more than 30 iterations are required to compute the  $j^{\text{th}}$  eigenvalue  $\lambda_j$ , then IERR is set to  $j$  and the routine terminates. In this case, if  $j < n$  then the eigenvalues  $\lambda_{j+1}, \dots, \lambda_n$  will have been computed and the results stored in the WR and WI arrays.

### Remarks.

- (1) Even though the balancing operation does not increase the theoretical bounds on the errors, nevertheless at times it may result in a slight loss of accuracy. On the other hand, balancing requires little additional time and in certain cases can improve the accuracy by as much as 5–6 significant digits. Thus it is recommended that balancing be done.
- (2) AR and AI are destroyed during computation. CEIG reduces the complex matrix (represented by AR and AI) to upper Hessenberg form with unitary similarity transformations. Then the QR algorithm is used to obtain the eigenvalues.

**Usage.** If one has a complex matrix  $A$ , then AR and AI can be obtained using the matrix subroutines CMREAL and CMIMAG.

**Programming.** CEIG is a driver routine for the EISPACK subroutines CBAL, CORTH, and COMQR. These subroutines were developed at Argonne National Laboratory. The functions SPMPAR and IPMPAR are also used.

**Reference.** Smith, B. T., Boyle, J. M., et al., *Matrix Eigensystem Routines – EISPACK Guide* (Second Edition), Springer-Verlag, 1976.



## COMPUTATION OF EIGENVALUES AND EIGENVECTORS OF COMPLEX MATRICES

The subroutine CEIGV is available for computing the eigenvalues and eigenvectors of complex matrices. This routine frequently yields values for the eigenvalues that are accurate to 13–14 significant digits. However, be aware that errors in the eigenvalues, no matter how seemingly insignificant, can be considerably magnified in the computation of the eigenvectors. It is not at all unusual to obtain an eigenvalue and eigenvector where the eigenvalue is correct to within 2–3 units of the 14<sup>th</sup> significant digit, but the components of the corresponding eigenvector are only accurate to 9–10 significant digits. If the eigenvalues of a matrix are not distinct or if they are exceedingly tightly clustered, then a severe drop in accuracy can occur. In this case one should not expect the eigenvalues to have more than 7–8 digit accuracy, and the situation regarding the eigenvectors is totally unpredictable. The components of such an eigenvector may be correct to 6–7 significant digits, or the eigenvector may not even be an eigenvector! In this case the results should be checked.

**CALL CEIGV**(IBAL,AR,AI,ka,n,WR,WI,ZR,ZI,IERR,TEMP)

AR and AI are real matrices of order  $n \geq 1$  and WR and WI are real arrays of dimension  $n$  or larger. AR and AI are the real and imaginary portions of the complex matrix whose eigenvalues and eigenvectors are to be computed. When CEIGV is called the eigenvalues  $\lambda_1, \dots, \lambda_n$  and corresponding eigenvectors  $z_1, \dots, z_n$  are computed. The real parts of the eigenvalues are stored in WR(1), ..., WR( $n$ ) and the imaginary parts are stored in WI(1), ..., WI( $n$ ). The eigenvalues are unordered.

It is assumed that the input argument  $ka$  is the number of rows in the dimension statements for AR and AI in the calling program. ZR and ZI are real arrays of dimension  $ka \times n$ . For  $j = 1, \dots, n$  the real parts of the components of the eigenvector  $z_j$  are stored in the  $j^{\text{th}}$  column of ZR (in locations ZR(1, $j$ ), ..., ZR( $n$ , $j$ )) and the imaginary parts are stored in the  $j^{\text{th}}$  column of ZI. The eigenvectors  $z_1, \dots, z_n$  are not normalized.

IBAL is an input argument that can be assigned any integer value. If IBAL  $\neq 0$  then the complex matrix (represented by AR and AI) is balanced before the eigenvalues and eigenvectors are computed. Otherwise, if IBAL = 0 then the complex matrix is not balanced.

TEMP is a real array used for temporary storage by the routine. If no balancing is to be done (i.e., if IBAL = 0) then TEMP must be of dimension  $2n$  or larger. Otherwise, if balancing is to be performed then TEMP must be of dimension  $3n$  or larger.

**Error Return.** IERR is an integer variable. If all the eigenvalues and eigenvectors are found then IERR is set to 0. Otherwise, if more than 30 iterations are required to compute the  $j^{\text{th}}$  eigenvalue  $\lambda_j$ , then IERR is set to  $j$  and the routine terminates. In this case, if  $j < n$  then the eigenvalues  $\lambda_{j+1}, \dots, \lambda_n$  will have been computed and the results stored in the WR and WI arrays. However, none of the eigenvectors will have been computed. The eigenvectors are computed only after all the eigenvalues have been obtained.

**Remarks.**

- (1) Even though the balancing operation does not increase the theoretical bounds on the errors, nevertheless at times it may result in a slight loss of accuracy. On the other hand, balancing requires little additional time and in certain cases can improve the accuracy of the eigenvalues by as much as 5-6 significant digits. When this occurs balancing will normally be needed to obtain the eigenvectors. In general, it is recommended that balancing be done.
- (2) AR and AI are destroyed during computation. CEIGV reduces the complex matrix (represented by AR and AI) to upper Hessenberg form with unitary similarity transformations. Then the QR algorithm is employed to obtain the eigenvalues, and back-substitution is performed to generate the eigenvectors.

**Usage.** If one has a complex matrix  $A$ , then AR and AI can be obtained using the matrix subroutines CMREAL and CMIMAG.

**Programming.** CEIGV is a driver routine for the EISPACK subroutines CBAL, CORTH, COMQR2, and CBABK2. These subroutines were developed at Argonne National Laboratory. The functions SPMPAR and IPMPAR are also used.

**Reference.** Smith, B. T., Boyle, J. M., et al., *Matrix Eigensystem Routines - EISPACK Guide* (Second Edition), Springer-Verlag, 1976.

## DOUBLE PRECISION COMPUTATION OF EIGENVALUES OF COMPLEX MATRICES

The subroutine DCEIG is available for the double precision computation of the eigenvalues of complex matrices. This routine frequently yields results accurate to 26–28 significant digits. However, if the eigenvalues are not distinct or if they are exceedingly tightly clustered, then a severe drop in accuracy can occur. In this case one should not expect more than 13–14 digit accuracy.

**CALL DCEIG**(IBAL,AR,AI,ka,n,WR,WI,IERR)

AR and AI are double precision matrices of order  $n \geq 1$ , and WR and WI are double precision arrays of dimension  $n$  or larger. AR and AI are the real and imaginary parts of the matrix whose eigenvalues are to be computed. When DCEIG is called the eigenvalues  $\lambda_1, \dots, \lambda_n$  are computed. The real parts of the eigenvalues are stored in WR(1), ..., WR( $n$ ) and the imaginary parts are stored in WI(1), ..., WI( $n$ ). The eigenvalues are unordered.

IBAL and ka are input arguments. It is assumed that ka is the number of rows in the dimension statements for AR and AI in the calling program. IBAL may be any integer. If IBAL  $\neq 0$  then the complex matrix (represented by AR and AI) is balanced before the eigenvalues are computed. Otherwise, if IBAL = 0 then the complex matrix is not balanced.

**Error Return.** IERR is an integer variable. If all the eigenvalues are found then IERR is set to 0. Otherwise, if more than 50 iterations are required to compute the  $j^{\text{th}}$  eigenvalue  $\lambda_j$ , then IERR is set to  $j$  and the routine terminates. In this case, if  $j < n$  then the eigenvalues  $\lambda_{j+1}, \dots, \lambda_n$  will have been computed and the results stored in the WR and WI arrays.

### Remarks.

- (1) AR and AI are destroyed during computation.
- (2) DCEIG is a double precision version of the eigenvalue routine CEIG.

**Programming.** DCEIG is a driver routine for the subroutines DCBAL, DCORTH, and DCOMQR. These subroutines are double precision versions of the EISPACK subroutines CBAL, CORTH, and COMQR, developed at Argonne National Laboratory. The double precision versions were prepared by A. H. Morris. The functions DCPABS, DPMPAR, IPMPAR and subroutine DCSQRT are also used.

**Reference.** Smith, B. T., Boyle, J. M., et al., *Matrix Eigensystem Routines – EISPACK Guide* (Second Edition), Springer-Verlag, 1976.

## DOUBLE PRECISION COMPUTATION OF EIGENVALUES AND EIGENVECTORS OF COMPLEX MATRICES

The subroutine DCEIGV is available for the double precision computation of the eigenvalues and eigenvectors of complex matrices. The routine frequently yields values for the eigenvalues that are accurate to 26–28 significant digits. However, be aware that the errors in the eigenvalues, no matter how seemingly insignificant, can be considerably magnified in the computation of the eigenvectors. If the eigenvalues are not distinct or if they are exceedingly tightly clustered, then a severe drop in accuracy can occur. In this case one should not expect the eigenvalues to have more than 13–14 digit accuracy.

**CALL DCEIGV**(IBAL,AR,AI,ka,n,WR,WI,ZR,ZI,IERR,TEMP)

AR and AI are double precision matrices of order  $n \geq 1$  and WR and WI are double precision arrays of dimension  $n$  or larger. AR and AI are the real and imaginary portions of the complex matrix whose eigenvalues and eigenvectors are to be computed. When DCEIGV is called the eigenvalues  $\lambda_1, \dots, \lambda_n$  and corresponding eigenvectors  $z_1, \dots, z_n$  are computed. The real parts of the eigenvalues are stored in WR(1), ..., WR( $n$ ) and the imaginary parts are stored in WI(1), ..., WI( $n$ ). The eigenvalues are unordered.

It is assumed that the input argument ka is the number of rows in the dimension statements for AR and AI in the calling program. ZR and ZI are double precision arrays of dimension  $ka \times n$ . For  $j = 1, \dots, n$  the real parts of the components of the eigenvector  $z_j$  are stored in the  $j^{\text{th}}$  column of ZR (in locations ZR(1, $j$ ), ..., ZR( $n$ , $j$ )) and the imaginary parts are stored in the  $j^{\text{th}}$  column of ZI. The eigenvectors  $z_1, \dots, z_n$  are not normalized.

IBAL is an input argument that can be assigned any integer value. If IBAL  $\neq 0$  then the complex matrix (represented by AR and AI) is balanced before the eigenvalues and eigenvectors are computed. Otherwise, if IBAL = 0 then the complex matrix is not balanced.

TEMP is a double precision array used for temporary storage by the routine. If no balancing is to be done (i.e., if IBAL = 0) then TEMP must be of dimension  $2n$  or larger. Otherwise, if balancing is to be performed then TEMP must be of dimension  $3n$  or larger.

**Error Return.** IERR is an integer variable. If all the eigenvalues and eigenvectors are found then IERR is set to 0. Otherwise, if more than 50 iterations are required to compute the  $j^{\text{th}}$  eigenvalue  $\lambda_j$ , then IERR is set to  $j$  and the routine terminates. In this case, if  $j < n$  then the eigenvalues  $\lambda_{j+1}, \dots, \lambda_n$  will have been computed and the results stored in the WR and WI arrays. However, none of the eigenvectors will have been computed. The eigenvectors are computed only after all the eigenvalues have been obtained.

### Remarks.

- (1) AR and AI are destroyed during computation.
- (2) DCEIGV is a double precision version of the eigenvalue/eigenvector routine CEIGV.

**Programming.** DCEIGV is a driver for the subroutines DCBAL, DCORTH, DCMQR2, and DCBABK. These subroutines are double precision versions of the EISPACK routines CBAL, CORTH, COMQR2, and CBABK2, developed at Argonne National Laboratory. The double precision versions were prepared by A. H. Morris. The functions DCPABS, DPMPAR, IPMPAR and subroutine DCSQRT are also used.

**Reference.** Smith, B. T., Boyle, J. M., et al., *Matrix Eigensystem Routines - EISPACK Guide* (Second Edition), Springer-Verlag, 1976.

## $\ell_1$ SOLUTION OF SYSTEMS OF LINEAR EQUATIONS WITH EQUALITY AND INEQUALITY CONSTRAINTS

Let  $A$  be a  $k \times n$  matrix,  $C$  an  $\ell \times n$  matrix, and  $E$  an  $m \times n$  matrix. Also let  $b, d$ , and  $f$  be column vectors of dimensions  $k, \ell$ , and  $m$  respectively. The following subroutine is available for obtaining a column vector  $x$  of dimension  $n$  which minimizes  $\|Ax - b\|_1 = \sum_{i=1}^k |A_i x - b_i|$  subject to the constraints

$$\begin{aligned} Cx &= d \\ Ex &\leq f. \end{aligned}$$

Here  $A_i$  denotes the  $i^{\text{th}}$  row of  $A$ , and  $Ex \leq f$  means that every component of  $Ex$  is less than or equal to the corresponding component of  $f$ .

**CALL CL1**( $k, \ell, m, n, Q, kq, \text{KODE}, \text{TOL}, \text{ITER}, X, \text{RES}, \text{RNORM}, \text{WK}, \text{IWK}$ )

It is assumed that  $k \geq 1, \ell \geq 0, m \geq 0$ , and  $n \geq 1$ .  $Q$  is a 2-dimensional array with  $kq$  rows and at least  $n + 2$  columns where  $kq \geq k + \ell + m + 2$ . The matrices  $A, C, E$  and vectors  $b, d, f$  are stored in the first  $k + \ell + m$  rows and  $n + 1$  columns of  $Q$  as follows:

$$Q = \begin{pmatrix} A & b \\ C & d \\ E & f \end{pmatrix}$$

$Q$  is modified by the routine.

**KODE** is a variable used for input/output purposes,  $X$  an array of dimension  $n + 2$  or larger, and **RES** an array of dimension  $k + \ell + m$  or larger. On input **KODE** is normally set by the user to 0. This indicates that  $\|Ax - b\|_1$  is to be minimized subject only to the constraints  $Cx = d$  and  $Ex \leq f$ . However, if it is also desired that one or more variables  $x_j$  satisfy  $x_j \leq 0$  or  $x_j \geq 0$ , or that one or more residuals  $b_i - A_i x$  satisfy  $b_i - A_i x \leq 0$  or  $b_i - A_i x \geq 0$ , then the user may set **KODE** to a nonzero value. If **KODE**  $\neq 0$  on input, then the user must also set  $X(j)$  and  $\text{RES}(i)$  to the values

$$\begin{aligned} X(j) &= \begin{cases} -1.0 & x_j \leq 0 \\ 0.0 & x_j \text{ is unrestricted} \\ 1.0 & x_j \geq 0 \end{cases} \\ \text{RES}(i) &= \begin{cases} -1.0 & b_i - A_i x \leq 0 \\ 0.0 & b_i - A_i x \text{ is unrestricted} \\ 1.0 & b_i - A_i x \geq 0 \end{cases} \end{aligned}$$

for  $j = 1, \dots, n$  and  $i = 1, \dots, k$  to indicate the additional constraints which are desired.

**RNORM** is a variable. When **CL1** is called, if a vector  $x$  is found that minimizes  $\|Ax - b\|_1$  subject to the desired constraints, then **KODE** = 0 on output and the solution  $x$

is stored in  $X$ . Also RNORM is assigned the value  $\|Ax - b\|_1$ ,  $b - Ax$  is stored in the first  $k$  locations of RES,  $d - Cx$  is stored in the next  $\ell$  locations, and  $f - Ex$  is stored in the last  $m$  locations.

When CL1 is called, a modified form of the simplex algorithm is used to minimize  $\|Ax - b\|_1$ . The arguments TOL and ITER control the use of this algorithm. The input argument TOL is a positive tolerance. CL1 will not pivot on any quantity whose magnitude is less than TOL. Normally the setting  $TOL = 10^{-2\nu/3}$  suffices where  $\nu$  is the number of decimal digits of accuracy available.

Frequently the routine requires less than  $5(k + \ell + m)$  iterations of the simplex algorithm to solve the problem. ITER is a variable used for input/output purposes. On input the user must set ITER to the maximum number of iterations that will be permitted. When the routine terminates, ITER has for its value the number of iterations that were performed.

On output CODE reports the status of the results. The routine assigns CODE one of the following values:

CODE = 0 The problem was solved.

CODE = 1 The problem has no solution.

CODE = 2 Sufficient accuracy cannot be maintained to solve the problem using the current value of TOL.

CODE = 3 The maximum number of iterations were performed. More iterations are needed.

When  $CODE \geq 1$  on output,  $X$  contains the last vector  $\bar{x}$  which was obtained,  $RNORM = \|A\bar{x} - b\|_1$ , and RES contains the vectors  $b - A\bar{x}$ ,  $d - C\bar{x}$ , and  $f - E\bar{x}$ .

WK is an array of dimension  $2(k + \ell + m + n)$  or larger, and IWK is an array of dimension  $3(k + \ell + m) + 2n$  or larger. WK and IWK are work spaces for the routine.

**Programming.** CL1 calls the subroutine KL1. CL1 was written by I. Barrodale and F. D. K. Roberts (University of Victoria, British Columbia, Canada).

#### References.

- (1) Barrodale, I. and Roberts, F. D. K., "An Improved Algorithm for Discrete  $\ell_1$  Linear Approximation," *SIAM J. Numer. Analysis* 10 (1973), pp. 839-848.
- (2) \_\_\_\_\_, "An Efficient Algorithm for Discrete  $\ell_1$  Linear Approximation with Linear Constraints," *SIAM J. Numer. Analysis* 15 (1978), pp. 603-611.
- (3) \_\_\_\_\_, "Algorithm 552, Solution of the Constrained  $\ell_1$  Linear Approximation Problem," *ACM Trans. Math Software* 6 (1980), pp. 231-235.

## LEAST SQUARES SOLUTION OF SYSTEMS OF LINEAR EQUATIONS

Given an  $m \times n$  matrix  $A$  and an  $m \times \ell$  matrix  $B$ . The column vectors  $b_1, \dots, b_\ell$  of  $B$  specify  $\ell$  distinct linear least squares problems

$$Ax_j = b_j \quad (j = 1, \dots, \ell).$$

This set of problems can be written in the form  $AX = B$  where  $X$  is the  $n \times \ell$  matrix having the column vectors  $x_1, \dots, x_\ell$ . There always exists a unique minimum length least squares solution  $x_j$  for each  $Ax_j = b_j$ . The subroutines LLSQ, HFTI, and HFTI2 are available for obtaining the minimum length solution matrix  $X$ . HFTI and HFTI2 are more general than LLSQ, being able to solve arbitrary systems  $AX = B$ . LLSQ assumes that  $m \geq n > 1$  and that the rank of  $A$  is  $n$ . The routines perform Householder triangularization. HFTI and HFTI2 require more time than LLSQ, but may be more accurate. In LLSQ all calculations are performed in single precision. In HFTI and HFTI2 most inner products are computed in double precision and the results stored in single precision.

**CALL LLSQ**( $m, n, A, ka, B, kb, \ell, WK, IWK, IERR$ )

It is assumed that  $m \geq n > 1$  and that the rank of  $A$  is  $n$ . The input arguments  $ka$  and  $kb$  have the following values:

$ka$  = the number of rows in the dimension statement for  $A$  in the calling program

$kb$  = the number of rows in the dimension statement for  $B$  in the calling program

It is required that  $ka \geq m$  and  $kb \geq m$ .

$IERR$  is an integer variable. When LLSQ is called, if no input errors are detected then  $IERR$  is set to 0 and the solution matrix  $X$  stored in  $B$ . Also, if  $m \neq n$  then the residual norm  $\|Ax_j - b_j\|$  is computed and stored in  $B(n+1, j)$  for  $j = 1, \dots, \ell$ .<sup>1</sup>

$WK$  and  $IWK$  are arrays of dimension  $n$  or larger that are work spaces for the routine.

**Error Return.**  $IERR \neq 0$  when  $m \geq n > 1$  is not satisfied ( $IERR = 1$ ) or the rank of  $A$  is less than  $n$  ( $IERR = 2$ ).

**Note.**  $A$  is destroyed during computation.

**Programming.** LLSQ is a driver for the subroutines ORTHO and ORSOL, written by Nai-Kuan Tsao and Paul J. Nikolai (Aerospace Research Laboratories, Wright-Patterson Air Force Base).

**Reference.** Tsao, N. K. and Nikolai, P. J., *Procedures using Orthogonal Transformations for Linear Least Squares Problems*. Report ARL TR 74-0124, Aerospace Research Laboratories, Wright-Patterson Air Force Base, 1974.

<sup>1</sup>Throughout this section  $\|c\| = \sqrt{\sum c_i^2}$  for any vector  $c = (c_1, \dots, c_m)$ .



**CALL HFTI**(*A, ka, m, n, B, kb, l, τ, k, RNORM, H, G, IP*)  
**CALL HFTI2**(*A, ka, m, n, B, kb, l, D, τ, k, RNORM, H, G, IP, IERR*)

It is assumed that  $m, n \geq 1$  and that the input arguments  $ka$  and  $kb$  have the following values:

$ka$  = the number of rows in the dimension statement for  $A$  in the calling program

$kb$  = the number of rows in the dimension statement for  $B$  in the calling program

It is required that  $ka \geq m$ . Also, if  $l \geq 1$  then  $kb \geq \max\{m, n\}$ .

If  $l \geq 1$  then  $RNORM$  is an array of dimension  $l$  or larger. When HFTI or HFTI2 is called, the minimum length solution matrix  $X$  is computed and stored in  $B$ . Also the residual norm  $\|Ax_j - b_j\|$  is computed and stored in  $RNORM(j)$  for  $j = 1, \dots, l$ .

$H, G$ , and  $IP$  are arrays of dimension  $n$  or larger that are work spaces for the routines.

### The parameters $\tau, k$ , and $D$ .

The argument  $\tau$  is a tolerance that is set by the user,  $k$  a variable, and  $D$  an array of dimension  $\min\{m, n\}$  or larger. It is assumed that  $\tau \geq 0$ . Normally  $\tau = 0$  is the setting that is used.  $D$  and  $k$  are set by the routines.

In order to understand the use of  $\tau, k$ , and  $D$  one must be briefly acquainted with the processing of  $A$ . The routines first reduce  $A$  to a triangular matrix  $C$  where  $A = QCP$ .  $Q$  is an orthogonal matrix and  $P$  a permutation matrix.  $P$  is defined so that the diagonal elements  $c_{ii}$  of  $C$  satisfy  $|c_{ii}| \geq |c_{i+1, i+1}|$  for each  $i$ . The variable  $k$  is set to the largest integer such that  $|c_{kk}| > \tau$ , and if HFTI2 is used then the diagonal elements  $c_{ii}$  are stored in  $D$ .  $C$  is now regarded as the partitioned matrix

$$C = \begin{pmatrix} C_1 & C_2 \\ 0 & C_3 \end{pmatrix}$$

where  $C_1$  is a  $k \times k$  matrix. Minimum length least squares solutions  $x_j$  are then computed for the problems  $Ax_j = b_j$  using only the first  $k$  rows of  $C$ . This is equivalent to replacing  $A$  with

$$\tilde{A} = Q \begin{pmatrix} C_1 & C_2 \\ 0 & 0 \end{pmatrix} P$$

and solving  $\tilde{A}x_j = b_j$  for  $j = 1, \dots, l$ .

Since  $|c_{11}| \geq \dots \geq |c_{kk}| > \tau$  clearly  $k$  is the rank of  $\tilde{A}$ . It is also true that the ratio  $|c_{11}|/|c_{kk}|$  is a lower bound on the condition number of  $C_1$  (relative to the spectral norm). Thus, if the ratio is extremely large (say  $\geq 10^8$ ) then a severe loss of accuracy can be expected. A large ratio may be due all or in part to rank deficiency (or near rank deficiency) of the matrix  $A$ . Fortunately, rank deficiency is frequently not too difficult to detect and cure. When  $A$  is rank deficient then machine roundoff may assign  $c_{kk}$  a small value, say  $10^{-14}$ , when it should be 0. The cure is to examine the diagonal elements  $c_{ii}$  which are stored in  $D$ , to reset  $\tau$  so as to eliminate the unwanted  $c_{ii}$ 's, and then to rerun the problem. This will reduce the order of  $C_1$ , thereby lowering the rank of the replacement matrix  $\tilde{A}$ .  $C_1$  will now be better conditioned, but the value of the residual norms  $\|Ax_j - b_j\|$  may be larger. If the norms do increase, then the solution obtained will be satisfactory only if the size of the increased norms fall within acceptable bounds.

**Remarks.**

- (1) The variable  $k$  is set to 0 if all  $|c_{ii}| \leq \tau$ . If  $k = 0$  then the zero matrix is the solution for  $AX = B$ .
- (2) If  $\ell \leq 0$  then the decomposition  $A = QCP$  is performed, the diagonal elements of  $C$  are stored in  $D$ , and  $k$  is computed.  $B$  and  $RNORM$  are ignored.
- (3) The contents of  $A$  are destroyed by the routines.
- (4) HFTI and HFTI2 yield the same results.

**Error Return.** IERR is a variable that is set by the routine. If no input errors are detected then IERR is set to 0. Otherwise, IERR is assigned one of the following values:

IERR = 1 if  $m > ka$

IERR = 2 if  $\ell > 1$  and  $kb < \max \{m, n\}$

When an error is detected, the routine immediately terminates.

**Programming.** HFTI and HFTI2 call the subroutine H12. These routines were written by Charles L. Lawson and Richard J. Hanson (Jet Propulsion Laboratory), and modified by A. H. Morris.

**Reference.** Lawson, C. L., and Hanson, R. J., *Solving Least Squares Problems*, Prentice-Hall, Inc., Englewood Cliffs, New Jersey, 1974.

## LEAST SQUARES SOLUTION OF OVERDETERMINED SYSTEMS OF LINEAR EQUATIONS WITH ITERATIVE IMPROVEMENT

Given an  $m \times n$  matrix  $A$  and an  $m \times \ell$  matrix  $B$ . The column vectors  $b_1, \dots, b_\ell$  of  $B$  specify  $\ell$  distinct linear least squares problems

$$Ax_j = b_j \quad (j = 1, \dots, \ell).$$

This set of problems can be written in the form  $AX = B$  where  $X$  is the  $n \times \ell$  matrix having the column vectors  $x_1, \dots, x_\ell$ . Assume that  $m \geq n > 1$  and that the rank of  $A$  is  $n$ . Then there exists a unique least squares solution  $x_j$  for each  $Ax_j = b_j$ . The subroutine LLSQMP is available for obtaining the solution matrix  $X$ . Iterative improvement is performed to compute  $X$  to machine accuracy.

**CALL LLSQMP**( $m, n, A, ka, B, kb, \ell, WK, IWK, IERR$ )

It is assumed that  $m \geq n > 1$  and that the rank of  $A$  is  $n$ . The input arguments  $ka$  and  $kb$  have the following values:

$ka$  = the number of rows in the dimension statement for  $A$  in the calling program

$kb$  = the number of rows in the dimension statement for  $B$  in the calling program

It is required that  $ka \geq m$  and  $kb \geq m$ .

When LLSQMP is called, the solution  $X$  is computed and stored in  $B$ . Also, if  $m \neq n$  then the residual norm  $\|Ax_j - b_j\|$  is computed and stored in  $B(n+1, j)$  for  $j = 1, \dots, \ell$ .<sup>1</sup>  $A$  is not modified by the routine.

$WK$  is an array of dimension  $mn + 2m + n$  or larger, and  $IWK$  an array of dimension  $n$  or larger.  $WK$  and  $IWK$  are work spaces for the routine.

$IERR$  is a variable that is set by the routine. When LLSQMP terminates,  $IERR$  has one of the following values:

$IERR = 0$  The solution  $X$  was computed to machine accuracy.

$IERR = 1$   $X$  was obtained, but not to machine accuracy.

$IERR = 2$  The restriction  $m \geq n > 1$  is not satisfied.

$IERR = 3$  The rank of  $A$  is less than  $n$ .

**Programming.** LLSQMP is a driver for the subroutines ORTHO, ORSOL, and ORIMP. These subroutines were written by Nai-Kuan Tsao and Paul J. Nikolai (Aerospace Research Laboratories, Wright-Patterson Air Force Base). ORIMP was modified by A. H. Morris. The function SPMPAR and subroutine MCOPI are also used.

**Reference.** Tsao, N. K. and Nikolai, P. J., *Procedures using Orthogonal Transformations for Linear Least Squares Problems*, Report ARL TR 74-0124, Aerospace Research Laboratories, Wright-Patterson Air Force Base, 1974.

<sup>1</sup> Here  $\|c\| = \sqrt{\sum_i c_i^2}$  for any vector  $c = (c_1, \dots, c_m)$ .

## DOUBLE PRECISION LEAST SQUARES SOLUTION OF SYSTEMS OF LINEAR EQUATIONS

Given an  $m \times n$  matrix  $A$  and an  $m \times \ell$  matrix  $B$ . The column vectors  $b_1, \dots, b_\ell$  of  $B$  specify  $\ell$  distinct linear least squares problems

$$Ax_j = b_j \quad (j = 1, \dots, \ell).$$

This set of problems can be written in the form  $AX = B$  where  $X$  is the  $n \times \ell$  matrix having the column vectors  $x_1, \dots, x_\ell$ . There always exists a unique minimum length least squares solution  $x_j$  for each  $Ax_j = b_j$ . The subroutines DLLSQ, DHFTI, and DHFTI2 are available for obtaining the minimum length solution matrix  $X$ . DHFTI and DHFTI2 are more general than DLLSQ, being able to solve arbitrary systems  $AX = B$ . DLLSQ assumes that  $m \geq n > 1$  and that the rank of  $A$  is  $n$ . The routines perform all calculations in double precision.

**CALL DLLSQ( $m, n, A, ka, B, kb, \ell, WK, IWK, IERR$ )**

$A$  and  $B$  are double precision arrays. It is assumed that  $m \geq n > 1$  and that the rank of  $A$  is  $n$ . The input arguments  $ka$  and  $kb$  have the following values:

$ka$  = the number of rows in the dimension statement for  $A$  in the calling program

$kb$  = the number of rows in the dimension statement for  $B$  in the calling program

It is required that  $ka \geq m$  and  $kb \geq m$ .

$IERR$  is an integer variable. When DLLSQ is called, if no input errors are detected then  $IERR$  is set to 0 and the solution matrix  $X$  stored in  $B$ . Also, if  $m \neq n$  then the residual norm  $\|Ax_j - b_j\|$  is computed and stored in  $B(n+1, j)$  for  $j = 1, \dots, \ell$ .<sup>1</sup>

$WK$  and  $IWK$  are arrays of dimension  $n$  or larger that are work spaces for the routine.  $WK$  is a double precision array.

**Error Return.**  $IERR \neq 0$  when  $m \geq n > 1$  is not satisfied ( $IERR = 1$ ) or the rank of  $A$  is less than  $n$  ( $IERR = 2$ ).

**Note.**  $A$  is destroyed during computation.

**Programming.** DLLSQ calls the subroutines DORTHO and DORSOL. These subroutines are double precision versions of ORTHO and ORSOL, written by Nai-Kuan Tsao and Paul J. Nikolai (Aerospace Research Laboratories, Wright-Patterson Air Force Base).

**Reference.** Tsao, N. K. and Nikolai, P. J., *Procedures using Orthogonal Transformations for Linear Least Squares Problems*. Report ARL TR 74-0124, Aerospace Research Laboratories, Wright-Patterson Air Force Base, 1974.

<sup>1</sup>Throughout this section  $\|c\| = \sqrt{\sum c_i^2}$  for any vector  $c = (c_1, \dots, c_m)$ .

**CALL DHFTI(A,ka,m,n,B,kb,l,r,k,RNORM,H,G,IP)**  
**CALL DHFTI2(A,ka,m,n,B,kb,l,D,r,k,RNORM,H,G,IP,IERR)**

$A$  and  $B$  are double precision arrays. It is assumed that  $m, n \geq 1$  and that the input arguments  $ka$  and  $kb$  have the following values:

$ka$  = the number of rows in the dimension statement for  $A$  in the calling program

$kb$  = the number of rows in the dimension statement for  $B$  in the calling program

It is required that  $ka \geq m$ . Also, if  $l \geq 1$  then  $kb \geq \max\{m, n\}$ .

If  $l \geq 1$  then  $RNORM$  is a double precision array of dimension  $l$  or larger. When  $DHFTI$  or  $DHFTI2$  is called, the minimum length solution matrix  $X$  is computed and stored in  $B$ . Also the residual norm  $\|Ax_j - b_j\|$  is computed and stored in  $RNORM(j)$  for  $j = 1, \dots, l$ .

$H, G$ , and  $IP$  are arrays of dimension  $n$  or larger that are work spaces for the routines.  $H$  and  $G$  are double precision arrays.

#### The parameters $r, k$ , and $D$ .

$r$  is a double precision tolerance that is set by the user,  $k$  an integer variable, and  $D$  a double precision array of dimension  $\min\{m, n\}$  or larger. It is assumed that  $r \geq 0$ . Normally  $r = 0$  is the setting that is used.  $D$  and  $k$  are set by the routines.

In order to understand the use of  $r, k$ , and  $D$  one must be briefly acquainted with the processing of  $A$ . The routines first reduce  $A$  to a triangular matrix  $C$  where  $A = QCP$ .  $Q$  is an orthogonal matrix and  $P$  a permutation matrix.  $P$  is defined so that the diagonal elements  $c_{ii}$  of  $C$  satisfy  $|c_{ii}| \geq |c_{i+1,i+1}|$  for each  $i$ . The variable  $k$  is set to the largest integer such that  $|c_{kk}| > r$ , and if  $DHFTI2$  is used then the diagonal elements  $c_{ii}$  are stored in  $D$ .  $C$  is now regarded as the partitioned matrix

$$C = \begin{pmatrix} C_1 & C_2 \\ 0 & C_3 \end{pmatrix}$$

where  $C_1$  is a  $k \times k$  matrix. Minimum length least squares solutions  $x_j$  are then computed for the problems  $Ax_j = b_j$  using only the first  $k$  rows of  $C$ . This is equivalent to replacing  $A$  with

$$\tilde{A} = Q \begin{pmatrix} C_1 & C_2 \\ 0 & 0 \end{pmatrix} P$$

and solving  $\tilde{A}x_j = b_j$  for  $j = 1, \dots, l$ .

Since  $|c_{11}| \geq \dots \geq |c_{kk}| > r$  clearly  $k$  is the rank of  $\tilde{A}$ . It is also true that the ratio  $|c_{11}| / |c_{kk}|$  is a lower bound on the condition number of  $C_1$  (relative to the spectral norm). Thus, if the ratio is extremely large (say  $\geq 10^8$ ) then a severe loss of accuracy can be expected. A large ratio may be due all or in part to rank deficiency (or near rank deficiency) of the matrix  $A$ . Fortunately, rank deficiency is frequently not too difficult to detect and cure. When  $A$  is rank deficient then machine roundoff may assign  $c_{kk}$  a small value, say  $10^{-28}$ , when it should be 0. The cure is to examine the diagonal elements  $c_{ii}$  which are stored in  $D$ , to reset  $r$  so as to eliminate the unwanted  $c_{ii}$ 's, and then to rerun the problem. This will reduce the order of  $C_1$ , thereby lowering the rank of the replacement

matrix  $\tilde{A}$ .  $C_1$  will now be better conditioned, but the value of the residual norms  $\|Ax_j - b_j\|$  may be larger. If the norms do increase, then the solution obtained will be satisfactory only if the size of the increased norms fall within acceptable bounds.

**Remarks.**

- (1) The variable  $k$  is set to 0 if all  $|c_{ii}| \leq \tau$ . If  $k = 0$  then the zero matrix is the solution for  $AX = B$ .
- (2) If  $\ell \leq 0$  then the decomposition  $A = QCP$  is performed, the diagonal elements of  $C$  are stored in  $D$ , and  $k$  is computed.  $B$  and  $RNORM$  are ignored.
- (3) The contents of  $A$  are destroyed by the routines.
- (4) DHFTI and DHFTI2 yield the same results. These routines are double precision versions of the subroutines HFTI and HFTI2.

**Error Return.** IERR is a variable that is set by the routine. If no input errors are detected then IERR is set to 0. Otherwise, IERR is assigned one of the following values:

IERR = 1 if  $m > ka$

IERR = 2 if  $\ell > 1$  and  $kb < \max \{m, n\}$

When an error is detected, the routine immediately terminates.

**Programming.** DHFTI and DHFTI2 call the subroutine DH12. These routines are modifications by A. H. Morris of the subroutines HFTI and H12, written by Charles L. Lawson and Richard J. Hanson (Jet Propulsion Laboratory).

**Reference.** Lawson, C. L., and Hanson, R. J., *Solving Least Squares Problems*, Prentice-Hall, Inc., Englewood Cliffs, New Jersey, 1974.

## LEAST SQUARES SOLUTION OF SYSTEMS OF LINEAR EQUATIONS WITH EQUALITY AND INEQUALITY CONSTRAINTS

Let  $A$  be an  $m_a \times n$  matrix,  $E$  an  $m_e \times n$  matrix,  $G$  an  $m_g \times n$  matrix,  $b$  a column vector of dimension  $m_a$ ,  $f$  a column vector of dimension  $m_e$ , and  $h$  a column vector of dimension  $m_g$ . The subroutine LSEI is available for finding a column vector  $x$  of dimension  $n$  that minimizes  $\|Ax - b\|$  subject to the constraints<sup>1</sup>

$$Ex \simeq f$$

$$Gx \geq h.$$

$Ex \simeq f$  states that  $x$  is a least squares solution of the equation  $Ex = f$ , and  $Gx \geq h$  means that every component of the vector  $Gx$  must be equal to or greater than the corresponding component of  $h$ . It is assumed that  $m_a \geq 0$ ,  $m_e \geq 0$ , and  $m_g \geq 0$ . If  $m_a = 0$  then LSEI solves  $Ex \simeq f$  subject to the constraints  $Gx \geq h$ .

**CALL LSEI**( $W, kw, m_e, m_a, m_g, n, OPT, x, RNORME, RNORMA,$   
 $IERR, WK, IWK$ )

If  $m = m_e + m_a + m_g$  then  $W$  is the  $m \times (n + 1)$  matrix:

$$W = \begin{pmatrix} E & f \\ A & b \\ G & h \end{pmatrix}$$

The input argument  $kw$  is assumed to have the value:

$kw =$  the number of rows in the dimension statement for  $W$  in the calling program

Thus it is required that  $kw \geq m$ .

$RNORME$  and  $RNORMA$  are real variables. When LSEI is called, if the constraints  $Ex \simeq f$  and  $Gx \geq h$  are consistent then  $x$  is computed,  $RNORME$  is assigned the value  $\|Ex - f\|$ , and  $RNORMA$  is assigned the value  $\|Ax - b\|$ .<sup>2</sup>

$OPT$  is an array, called the *option vector*, which permits the user to take advantage of certain options that are supplied by the routine. If no options are desired then  $OPT$  may be declared to have dimension 1 and  $OPT(1)$  must be assigned the value 1. The details concerning the available options and how to specify them in  $OPT$  are given below.

$IWK$  is an array of dimension  $m_g + 2n + 2$  or larger, and  $WK$  is an array of dimension  $2(m_e + n) + \max\{m_a + m_g, n\} + (m_g + 2)(n + 7)$  or larger.  $IWK$  and  $WK$  are work spaces. When LSEI is called, using a solution for  $Ex \simeq f$ , a reduced least squares problem with inequality constraints is obtained and solved. When the routine terminates  $IWK(1)$ ,  $IWK(2)$ ,  $IWK(3)$  contain the following information:

<sup>1</sup>Throughout this section  $\|c\|$  denotes the norm  $\sqrt{\sum_i c_i^2}$  for any vector  $c = (c_1, c_2, \dots)$ .

<sup>2</sup>If  $m_e = 0$  then  $RNORME = 0$ , and if  $m_a = 0$  then  $RNORMA = 0$ .

IWK(1) = the estimated rank of the matrix  $E$   
 IWK(2) = the estimated rank of the reduced problem  
 IWK(3) = the amount of storage in the array WK that was actually needed

IERR is a variable that is set by the routine. When LSEI terminates, IERR has one of the following values:

IERR = 0 The solution  $x$  was obtained. The equality  $Ex = f$  is satisfied when  $m_e \neq 0$ .  
 IERR = 1 The solution  $x$  was obtained. In this case  $\|Ex - f\| > 0$ .  
 IERR = 2 The problem cannot be solved. The constraints are inconsistent.  
 IERR = 4 (Input error) Either  $kw < m$ , the covariance matrix is requested and  $kw < n$ , or the option vector OPT is not defined properly.

If  $IERR \geq 2$  then  $x$ , RNORME, and RNORMA are not defined.

#### Remarks.

- (1)  $W$  is modified by the routine.
- (2) If  $m_e + m_a \leq 0$  or  $n \leq 0$  then IERR is set to 0 and the routine terminates.

**The option vector OPT.** If no options are desired then OPT may be declared to be of dimension 1 and OPT(1) must have the value 1. Otherwise, OPT is a linked list consisting of groups of data  $link_i$ ,  $key_i$ ,  $data_i$  ( $i = 1, \dots, s$ ). Each  $link_i$  and  $key_i$  is an integer. The amount of storage required by  $data_i$  depends on the value of  $key_i$ . The general layout of OPT is as follows:

OPT(1) =  $link_1$  (index of the first entry of the next group)  
 OPT(2) =  $key_1$  (key to the option)  
 OPT(3) = the first word of the data ( $data_1$ ) for this option  
 .  
 .  
 .  
 OPT( $link_1$ ) =  $link_2$  (index of the first entry of the next group)  
 OPT( $link_1 + 1$ ) =  $key_2$  (key to the option)  
 OPT( $link_1 + 2$ ) = the first word of the data ( $data_2$ ) for this option  
 .  
 .  
 .  
 OPT( $link_s$ ) = 1.0 (There are no more options to be considered.)

The following options are available:

key = 1 It is assumed that  $kw \geq n$ . Compute the  $n \times n$  covariance matrix and store it in the first  $n$  rows and columns of  $W$ . The data for this option is a single value. It must be nonzero for the covariance matrix to be computed.  
 key = 2 Scale the nonzero columns of the matrix  $\begin{pmatrix} E \\ A \\ G \end{pmatrix}$  so that they have



length 1. The data for this option is a single value. It must be nonzero for the scaling to be performed.

key = 3 Scale the columns of the matrix  $\begin{pmatrix} E \\ A \\ G \end{pmatrix}$ . The data for this option consists of  $n$  scaling factors, one for each matrix column.

key = 4 Change the internal tolerance  $\tau$  which is used for determining the rank of  $E$ . The data for this option is the new tolerance.  $\tau$  may be set to any value  $\geq \epsilon$  where  $\epsilon$  is the smallest floating point number for which  $1 + \epsilon > 1$  ( $\epsilon = 2^{-47}$  for the CDC 6700). If the new value is less than  $\epsilon$  then it is ignored and  $\tau$  is set to  $\epsilon$ . The default value employed for  $\tau$  is  $\sqrt{\epsilon}$ .

key = 5 Change the internal tolerance  $\tau$  which is used for rank determination in the reduced least squares problem. The data for this option is the new tolerance.  $\tau$  may be set to any value  $\geq \epsilon$  where  $\epsilon$  is the smallest floating point number for which  $1 + \epsilon > 1$ . If the new value is less than  $\epsilon$  then it is ignored and  $\tau$  is set to  $\epsilon$ . The default value employed for  $\tau$  is  $\sqrt{\epsilon}$ .

Also the key 8 and 9 options for the least squares subroutine WNNLS are permitted. (WNNLS is employed by LSEI.) The order of the options in the array OPT is arbitrary. If an option has an unrecognized key then the option is ignored. It is assumed that the dimension of OPT is no greater than 100000 and that the number of options is  $\leq 1000$ . If either of these assumptions is violated then IERR is set to 4 and the routine terminates. It is also required that  $\text{link}_i \neq \text{link}_j$  for  $i \neq j$ . If this restriction is not satisfied then the linked list OPT is circular and we again have an IERR = 4 error.

**Example.** Assume that we have an array  $D$  containing  $n$  scaling factors for the columns of the matrix  $\begin{pmatrix} E \\ A \\ G \end{pmatrix}$ , and that the tolerance TOL is always to be used for rank determination. Then OPT will have to be of dimension  $\geq n + 9$  and OPT can be defined as follows:

```

OPT(1) = N + 3           (Scaling option)
OPT(2) = 3.0
DO 10 I = 1, N
10  OPT(I + 2) = D(I)
OPT(N + 3) = N + 6       (Matrix E tolerance option)
OPT(N + 4) = 4.0
OPT(N + 5) = TOL
OPT(N + 6) = N + 9       (Reduced problem tolerance option)
OPT(N + 7) = 5.0
OPT(N + 8) = TOL
OPT(N + 9) = 1.0         (There are no more options.)

```

#### Remarks.

- (1) LSEI may perform poorly if the norms of the rows of  $A$  and  $E$  differ by many orders of magnitude, or if the norms of the rows of  $E$  are exceedingly small.
- (2) The covariance matrix obtained by the key = 1 option may not be meaningful when

there are inequality constraints  $Gx \geq h$ . This matrix assumes that any inequalities which are selected by the algorithm to be equalities remain equalities when the solution is perturbed. This, of course, may not be the case.

**Programming.** LSEI employs the subroutine LSI, LPDP, WNNLS, WNLSM, and WNLIT. These routines were written by Karen H. Haskell and Richard J. Hanson (Sandia Laboratories) and modified by A. H. Morris. The subroutines HFTI, H12, SROTM, SROTMG, SCOPY, SSWAP, SSCAL, SAXPY and functions SPMPAR, SDOT, SASUM, SNRM2, ISAMAX are also used.

#### References.

- (1) Hanson, R. J. and Haskell, K. H., "Algorithm 587: Two Algorithms for the Linearly Constrained Least Squares Problem," *ACM Trans. Math Software* 8 (1982), pp. 323-333.
- (2) Haskell, K. H. and Hanson, R. J., "An Algorithm for Linear Least Squares Problems with Equality and Nonnegativity Constraints," *Math. Programming* 21 (1981), pp. 98-118.

## LEAST SQUARES SOLUTION OF SYSTEMS OF LINEAR EQUATIONS WITH EQUALITY AND NONNEGATIVITY CONSTRAINTS

Let  $A$  be an  $m_a \times n$  matrix,  $E$  an  $m_e \times n$  matrix,  $b$  a column vector of dimension  $m_a$ , and  $f$  a column vector of dimension  $m_e$ . The subroutine WNNLS is available for finding a column vector  $x = (x_1, \dots, x_n)^t$  that minimizes  $\|Ax - b\|$  subject to the constraints<sup>1</sup>

$$\begin{aligned} Ex &\simeq f \\ x_i &\geq 0 \text{ for } i > \ell. \end{aligned}$$

$Ex \simeq f$  states that  $x$  is a least squares solution of the equation  $Ex = f$ . It is assumed that  $m_a \geq 0, m_e \geq 0$ , and  $0 \leq \ell \leq n$ . If  $m_a = 0$  then WNNLS solves  $Ex \simeq f$  subject to the constraints  $x_i \geq 0$  ( $i > \ell$ ).

**CALL WNNLS**( $W, kw, m_e, m_a, n, \ell, \text{OPT}, x, \text{RNORM}, \text{MODE}, \text{IWK}, \text{WK}$ )

If  $m = m_e + m_a$  then  $W$  is the  $m \times (n + 1)$  matrix:

$$W = \begin{pmatrix} E & f \\ A & b \end{pmatrix}$$

The input argument  $kw$  is assumed to have the value:

$kw$  = the number of rows in the dimension statement for  $W$  in the calling program  
Thus it is required that  $kw \geq m$ .

RNORM is a variable. When WNNLS is called,  $x$  is computed and RNORM is assigned the value  $\sqrt{\|Ax - b\|^2 + \|Ex - f\|^2}$ .<sup>2</sup>

OPT is an array, called the *option vector*, which permits the user to take advantage of certain options that are supplied by the routine. If no options are desired then OPT may be declared to have dimension 1 and OPT(1) must be assigned the value 1. The details concerning the available options and how to specify them in OPT are given below.

IWK is an array of dimension  $m + n$  or larger, and WK is an array of dimension  $m + 5n$  or larger. IWK and WK are work spaces for the routine.

**Error Return.** MODE is an integer variable that is set by the routine. If the problem is solved then MODE is assigned the value 0. Otherwise, MODE is assigned one of the following values:

MODE = 1 The maximum number of iterations ( $3(n - \ell)$  iterations) was exceeded. An approximate solution and its residual norm are stored in  $x$  and RNORM.

MODE = 2 (Input error) Either  $kw \geq m$  or  $0 \leq \ell \leq n$  is violated, or the option vector OPT is not properly defined.

<sup>1</sup>Throughout this section  $\|c\|$  denotes the norm  $\sqrt{\sum c_i^2}$  for any vector  $c = (c_1, c_2, \dots)$ .

<sup>2</sup>If  $m_a = 0$  then  $\text{RNORM} = \|Ex - f\|$ , and if  $m_e = 0$  then  $\text{RNORM} = \|Ax - b\|$ .

When an input error is detected, the routine immediately terminates. In this case  $x$  and RNORM are not defined.

**Remarks.**

- (1)  $W$  is modified by the routine.
- (2) If  $m \leq 0$  or  $n \leq 0$  then MODE is set to 0 and the routine terminates.

**The option vector OPT.** If no options are desired then OPT may be declared to be of dimension 1 and OPT(1) must have the value 1. Otherwise, OPT is a linked list consisting of groups of data  $link_i$ ,  $key_i$ ,  $data_i$  ( $i = 1, \dots, s$ ). Each  $link_i$  and  $key_i$  is an integer. The amount of storage required by  $data_i$  depends on the value of  $key_i$ . The general layout of OPT is as follows:

```

OPT(1) = link1      (index of the first entry of the next group)
OPT(2) = key1       (key to the option)
OPT(3) = the first word of the data (data1) for this option
.
.
.
OPT(link1) = link2   (index of the first entry of the next group)
OPT(link1 + 1) = key2 (key to the option)
OPT(link1 + 2) = the first word of the data (data2) for this option
.
.
.
OPT(links) = 1.0      (There are no more options to be considered.)

```

The following options are permitted:

- key = 6 Scale the nonzero columns of the matrix  $\begin{pmatrix} E \\ A \end{pmatrix}$  so that they have length 1. The data for this option is a single value. It must be nonzero for the scaling to be performed.
- key = 7 Scale the columns of the matrix  $\begin{pmatrix} E \\ A \end{pmatrix}$ . The data for this option consists of  $n$  scaling factors, one for each matrix column.
- key = 8 Change the internal tolerance  $\tau$  which is used for rank determination. The data for this option is the new tolerance.  $\tau$  may be set to any value  $\geq \epsilon$  where  $\epsilon$  is the smallest floating point number for which  $1 + \epsilon > 1$ . ( $\epsilon = 2^{-47}$  for the CDC 6700.) If the new value is less than  $\epsilon$  then it is ignored and  $\tau$  is set to  $\epsilon$ . The default value employed for  $\tau$  is  $\sqrt{\epsilon}$ .
- key = 9 Change the parameter BLOWUP. The reciprocal of this parameter is used in determining when solution components are too large. The data for this option is the new value for BLOWUP. It is assumed that BLOWUP < 1. BLOWUP may be set to any value  $\geq \epsilon$  where  $\epsilon$  is the smallest number for which  $1 + \epsilon > 1$ . If the new value is less than  $\epsilon$  then it is ignored and BLOWUP is set to  $\epsilon$ . The default value used for BLOWUP is  $\sqrt{\epsilon}$ .

The order of the options in the array OPT is arbitrary. If an option has an unrecognized key then the option is ignored. It is assumed that the dimension of OPT is no greater than 100000 and that the number of options  $\leq 1000$ . If either of these assumptions is violated then MODE is set to 2 and the routine terminates. It is also required that  $\text{link}_i \neq \text{link}_j$  for  $i \neq j$ . If this restriction is not satisfied then the linked list OPT is circular and we again have a MODE = 2 error.

**Example.** Assume that we have an array  $D$  containing  $n$  scaling factors for the columns of the matrix  $\begin{pmatrix} E \\ A \end{pmatrix}$ , and that TOL is the tolerance to be used for rank determination. Then OPT will have to be of dimension  $\geq n + 6$  and OPT can be defined as follows:

```

      OPT(1) = N + 3           (Scaling option)
      OPT(2) = 7.0
      DO 10 I = 1, N
10    OPT(I + 2) = D(I)
      OPT(N + 3) = N + 6      (Tolerance option)
      OPT(N + 4) = 8.0
      OPT(N + 5) = TOL
      OPT(N + 6) = 1.0       (There are no more options.)

```

**Remark.** WNNLS may perform poorly if the norms of the rows of  $A$  and  $E$  differ by many orders of magnitude, or if the norms of the rows of  $E$  are exceedingly small.

**Programming.** WNNLS employs the subroutines WNLSM and WNLIT. These routines were written by Karen H. Haskell and Richard J. Hanson (Sandia Laboratories), and modified by A. H. Morris and Virgis Dadurkevicius (Astronomical Observatory, Vilnius University, Lithuania). The subroutines H12, SROTM, SROTMG, SCOPY, SSWAP, SSCAL, SAXPY and functions SPMPAR, SASUM, SNRM2, ISAMAX are also used.

#### References.

- (1) Dadurkevicius, V., "Remark on Algorithm 587," *ACM Trans. Math Software* 15 (1989), pp. 257-261.
- (2) Hanson, R. J. and Haskell, K. H., "Algorithm 587: Two Algorithms for the Linearly Constrained Least Squares Problem," *ACM Trans. Math Software* 8 (1982), pp. 323-333.
- (3) Haskell, K. H. and Hanson, R. J., "An Algorithm for Linear Least Squares Problems with Equality and Nonnegativity Constraints," *Math. Programming* 21 (1981), pp. 98-118.

## LEAST SQUARES ITERATIVE IMPROVEMENT SOLUTION OF SYSTEMS OF LINEAR EQUATIONS WITH EQUALITY CONSTRAINTS

Let  $A$  be an  $m_a \times n$  matrix,  $E$  an  $m_e \times n$  matrix,  $B$  an  $m_a \times \ell$  matrix, and  $F$  an  $m_e \times \ell$  matrix. It is assumed that  $0 \leq m_e \leq n$  and that the rank of  $E$  is  $m_e$ . Let  $b_1, \dots, b_\ell$  denote the column vectors of  $B$  and  $f_1, \dots, f_\ell$  the column vectors of  $F$ . The subroutine L2SLV is available for finding the unique minimum length column vector  $x_j$  of dimension  $n$  that minimizes  $\|Ax_j - b_j\|$  subject to the equality constraints  $Ex_j = f_j$  (if there are any) for  $j = 1, \dots, \ell$ .<sup>1</sup> Iterative improvement is performed to compute the vectors  $x_1, \dots, x_\ell$  to machine accuracy. It is assumed that  $m_a \geq 0$ . If  $m_a = 0$  then L2SLV finds the unique minimum length solution  $x_j$  to  $Ex_j = f_j$  for  $j = 1, \dots, \ell$ .

CALL L2SLV( $m, n, m_e, \ell, \tilde{A}, ka, \tilde{B}, kb, \text{WGTS}, \text{TOL}, N1, \text{IPIVOT},$   
 $X, kx, R, kr, T, kt, \text{WK}, \text{IERR}$ )

If  $m = m_e + m_a$  then  $\tilde{A}$  is the  $m \times n$  matrix  $\begin{pmatrix} E \\ A \end{pmatrix}$  and  $\tilde{B}$  is the  $m \times \ell$  matrix  $\begin{pmatrix} F \\ B \end{pmatrix}$ . The input arguments  $ka$  and  $kb$  have the values:

$ka$  = the number of rows in the dimension statement for  $\tilde{A}$  in the calling program

$kb$  = the number of rows in the dimension statement for  $\tilde{B}$  in the calling program

It is assumed that  $m \geq 1, n \geq 1, \ell \geq 1, ka \geq m$ , and  $kb \geq m$ .  $\tilde{A}$  and  $\tilde{B}$  are not modified by the routine.

WGTS is an array containing  $m$  nonnegative weights. The first  $m_e$  weights are set to 1.0 by the routine. Let  $w_1, \dots, w_{m_a}$  denote the remaining weights (i.e., let  $w_i = \text{WGTS}(m_e + i)$  for  $i = 1, \dots, m_a$ ). The remaining weights are supplied by the user. In effect,  $w_i$  is the weighting that is given to the  $i^{\text{th}}$  equation in the least squares problem  $Ax_j = b_j$ . If  $W$  denotes the  $m_a \times m_a$  diagonal matrix  $\text{diag}(w_1, \dots, w_{m_a})$  then L2SLV finds the unique minimum length vector that minimizes  $\|WAx_j - Wb_j\|$  subject to  $Ex_j = f_j$  for  $j = 1, \dots, \ell$ . For convenience,  $\tilde{W}$  will denote the  $m \times m$  diagonal matrix  $\text{diag}(1, \dots, 1, w_1, \dots, w_{m_a})$ .

$X$  is an  $n \times \ell$  matrix that contains the solution vectors  $x_1, \dots, x_\ell$  when the routine terminates. The input argument  $kx$  is the number of rows in the dimension statement for  $X$  in the calling program. It is assumed that  $kx \geq n$ .

$R$  is an  $m \times \ell$  matrix. Let  $\tilde{b}_j$  denote the  $j^{\text{th}}$  column vector  $\begin{pmatrix} f_j \\ b_j \end{pmatrix}$  of  $\tilde{B}$  for  $j = 1, \dots, \ell$ . Then L2SLV stores the residual vector  $r_j = \tilde{W}\tilde{b}_j - \tilde{W}\tilde{A}x_j$  in the  $j^{\text{th}}$  column of  $R$ . The input argument  $kr$  is the number of rows in the dimension statement for  $R$  in the calling program. It is assumed that  $kr \geq m$ .

WK is an array of dimension  $6(m + n) + 2\ell$  or larger that is used for a work space. When L2SLV terminates, for  $j = 1, \dots, \ell$

$$WK(j) = \begin{cases} n_j & \text{if iterative improvement of the solution } x_j \text{ converged} \\ -n_j & \text{if iterative improvement of } x_j \text{ failed to converge} \end{cases}$$

<sup>1</sup>Throughout this section  $\|c\|$  denotes the norm  $\sqrt{\sum_i c_i^2}$  for any vector  $c = (c_1, \dots, c_{m_a})$ .

where  $n_j$  is the number of iterations in the iterative improvement process that were performed in computing  $x_j$ . Also

$WK(\ell+j)$  = the estimated number of correct digits in  $x_j$  before iterative improvement was performed  
for  $j = 1, \dots, \ell$ .

TOL and N1 correspond to the parameters  $\tau$  and  $k$  in the least squares subroutines HFTI and HFTI2. TOL is a nonnegative number that is specified by the user, and N1 is a variable that is set by the routine. When L2SLV is called, modified Gram-Schmidt orthogonalization with column pivoting is used to reduce  $\tilde{W}\tilde{A}$  to the form  $(A_1 A_2)$  where  $A_1$  is an  $m \times N_1$  matrix having rank  $N_1$ .  $A_1$  is of the form  $QU$  where  $Q^t Q = \text{diag}(d_1, \dots, d_{N_1})$  and  $U$  is an upper unit triangular matrix. The values  $d_1, d_2, \dots$  correspond to the diagonal elements  $c_{11}, c_{22}, \dots$  generated by HFTI and HFTI2 ( $d_i = c_{ii}^2$  for  $i = 1, 2, \dots$ ). The values are ordered so that  $d_i \geq d_{i+1}$ , and  $d_1, d_2, \dots$  are stored in  $WK(2\ell + 1)$ ,  $WK(2\ell + 2)$ , .... If  $m = m_e$  then N1 is assigned the value  $m_e$ . Otherwise, if  $m > m_e$  then N1 is the largest integer  $k$  for which  $d_k > \tau$ . Here  $\tau = \text{TOL}$  if  $\text{TOL} > 0$ , and  $\tau = (n\epsilon)^2 d_{m_e+1}$  where  $\epsilon$  is the smallest value for which  $1 + \epsilon > 1$  ( $\epsilon = 2^{-47}$  on the CDC 6700) if  $\text{TOL} = 0$ . Thus, if  $\text{TOL} = 0$  then a tolerance based on the computer precision is used to determine the rank of  $N_1$  of  $\tilde{W}\tilde{A}$ . Otherwise, if  $\text{TOL} > 0$  then TOL is the tolerance that is used to specify the rank of the problem to be solved. If the user inadvertently sets TOL to be negative then L2SLV resets TOL to be 0.

IPIVOT is an array of dimension  $n$  or larger that is used by L2SLV to record the order in which the columns of  $\tilde{W}\tilde{A}$  are selected by the pivoting procedure when  $\tilde{W}\tilde{A}$  is reduced to  $(A_1 A_2)$ . If  $N_1 < n$  then the first  $N_1$  elements of IPIVOT are the indices of the columns of  $\tilde{W}\tilde{A}$  from which the matrix  $A_1$  is generated.

$T$  is a 2-dimensional array of dimension  $kt \times n$  that is used for temporary storage. It is assumed that  $kt \geq m + n$ . When L2SLV terminates, if  $N_1 = n$  then the unscaled  $n \times n$  covariance matrix is stored in the first  $n$  rows and columns of  $T$ . Iterative improvement is not performed on the covariance matrix.

**Error Return.** IERR is an integer variable that is set by the routine. If no input errors are detected and the results appear to be satisfactory, then IERR is set to 0. Otherwise, IERR is assigned one of the following values:

- IERR = 1 Either  $m, n$ , or  $\ell$  is not positive.
- IERR = 2 The restriction  $0 \leq m_e \leq \min\{m, n\}$  is not satisfied.
- IERR = 3 Either  $ka \geq m, kb \geq m, kx \geq n, kr \geq m$ , or  $kt \geq m + n$  is violated.
- IERR = 4  $WGTS(i)$  is negative for some  $i > m_e$ .
- IERR = 5 Either  $\tilde{W}\tilde{A} = 0$  or  $E = 0$ .
- IERR = 6 The rank of  $E$  is less than  $m_e$ .
- IERR = 7 Iterative improvement of all the solutions  $x_1, \dots, x_\ell$  failed to converge.
- IERR = 8 Iterative improvement of one or more solutions failed to converge.
- IERR = 9 More than  $\mu$  iterations of the iterative improvement procedure were performed in computing some  $x_j$ . (Here it is assumed that a

$\mu$  decimal digit floating-point arithmetic is being used.  $\mu = 14$  for the CDC 6700.)

- IERR = 10 The accuracy of some  $x_j$  before iterative improvement was estimated to be less than half a decimal digit.
- IERR = 11 One or more of the computed diagonal elements of the covariance matrix is negative. This is due to roundoff error. Theoretically, all the diagonal elements should be nonnegative. No evidence of severe ill-conditioning was detected.
- IERR = 12 One or more of the computed diagonal elements of the covariance matrix is negative. This is due to roundoff error. Theoretically, all the diagonal elements should be nonnegative. The problem appears to be extremely ill-conditioned.

When an input error is detected (IERR = 1, 2, ..., 6) then L2SLV immediately terminates. If evidence of severe ill-conditioning is detected, then IERR is set to 8, 9, or 10 and computation of the solutions continues. If iterative improvement appears to converge for one or more of the solutions, then the covariance matrix is also computed (when  $N1 = n$ ). However, if iterative improvement fails for all the solutions  $x_1, \dots, x_\ell$  then IERR is set to 7 and the covariance matrix is not computed.

**Note.** WK(1), ..., WK(2 $\ell$ ) should be examined when severe ill-conditioning is detected.

**Programming.** L2SLV employs the subroutines DECOM2, SOLVE2, SOLVE3, and COVAR. These routines were written by Roy Wampller (National Bureau of Standards). L2SLV is a slightly modified version by A. H. Morris of the subroutine L2B discussed in reference (4). The algorithm employed for finding and iteratively improving the least squares solutions is described in references (1)–(3). The function SPMPAR is also used.

#### References.

- (1) Bjorck, Ake, "Solving Linear Least Squares Problems by Gram-Schmidt Orthogonalization," *BIT* **7** (1967), pp. 1–21.
- (2) \_\_\_\_\_, "Iterative Refinement of Linear Least Squares Solutions I," *BIT* **7** (1967), pp. 257–278.
- (3) \_\_\_\_\_, "Iterative Refinement of Linear Least Squares Solutions II," *BIT* **8** (1968), pp. 8–30.
- (4) Wampller, Roy, "Solutions to Weighted Least Squares Problems by Modified Gram-Schmidt with Iterative Refinement," *ACM Trans. Math Software* **5** (1979), pp. 457–465.



## ITERATIVE LEAST SQUARES SOLUTION OF BANDED LINEAR EQUATIONS

Given an  $m \times n$  matrix  $A$ , a column vector  $b$  of dimension  $m$ , and a real number  $\lambda$ . Let  $\bar{A} = \begin{pmatrix} A \\ \lambda I \end{pmatrix}$  where  $I$  is the  $n \times n$  identity matrix, and let  $\bar{b} = \begin{pmatrix} b \\ 0 \end{pmatrix}$ . The problem is to find a column vector  $x$  of dimension  $n$  which is a least squares solution of  $\bar{A}x = \bar{b}$ . If  $A$  is stored in band form then the following subroutine is available for solving this problem.

**CALL BLSQ**( $m, n, A, ka, m_l, m_u, \lambda, b, x, ATOL, BTOL, CONLIM, MXITER, IND, ITER, COND, RNORM, XNORM, WK$ )

$A$  is an  $m \times n$  matrix stored in band form,  $m_l$  the number of diagonals below the main diagonal containing nonzero elements, and  $m_u$  the number of diagonals above the main diagonal containing nonzero elements. The argument  $ka$  is the number of rows in the dimension statement for  $A$  in the calling program. It is assumed that  $0 \leq m_l < m, 0 \leq m_u < n$ , and  $ka \geq m$ . When BLSQ is called, an iterative procedure is used to obtain a least squares solution  $x$  of  $\bar{A}x = \bar{b}$ . The vector  $b$  is modified by the routine.

ATOL and BTOL are input arguments which specify the relative accuracy of  $A$  and  $b$  respectively. For example, if it is estimated that  $b$  is accurate to  $k$  decimal digits then one may set  $BTOL = 10^{-k}$ . It is required that  $ATOL \geq 0$  and  $BTOL \geq 0$ . If  $ATOL = 0$  or  $BTOL = 0$ , then it is assumed that  $A$  or  $b$  is accurate to machine precision.

Let  $\text{cond}(\bar{A})$  denote the condition number of  $\bar{A}$  relative to the Frobenius norm.<sup>1</sup> In each iteration of the algorithm being used, an estimate is made of the condition number  $\text{cond}(\bar{A})$ . The estimates form a monotonically nondecreasing sequence. The input argument CONLIM is an upper limit on  $\text{cond}(\bar{A})$ . If  $CONLIM > 0$  then BLSQ terminates when an estimate of  $\text{cond}(\bar{A})$  exceeds CONLIM. This termination may be needed to prevent small or zero singular values of  $A$  from coming into effect and causing damage to the solution  $x$ . CONLIM may be ignored by being set to 0. It is assumed that  $CONLIM \geq 0$ .

The input argument MXITER is the maximum number of iterations that are permitted. Normally BLSQ requires less than  $4n$  iterations. The related argument ITER is a variable. When the routine terminates  $ITER =$  the number of iterations that were performed.

COND, RNORM, and XNORM are variables. When BLSQ terminates  $COND =$  the last estimate made for  $\text{cond}(\bar{A})$ ,  $RNORM = \|\bar{A}x - \bar{b}\|$ , and  $XNORM = \|x\|$ .<sup>2</sup>

WK is an array of dimension  $2n$  or larger that is a work space for the routine.

The equations  $\bar{A}x = \bar{b}$  are considered to be *compatible* if for any least squares solution  $x$ ,  $\|\bar{A}x - \bar{b}\| = 0$ . IND is a variable that reports the status of the results. When BLSQ terminates, IND has one of the following values:

<sup>1</sup> $\text{cond}(\bar{A}) = \|\bar{A}\|_F \|\bar{A}^+\|_F$  where  $\bar{A}^+$  is the pseudoinverse of  $\bar{A}$ . Here  $\|C\|_F = \sqrt{\sum c_{ij}^2}$  for any matrix  $C = (c_{ij})$ .

<sup>2</sup> $\|c\| = \sqrt{\sum c_i^2}$  for any vector  $c = (c_1, c_2, \dots)$ .

- IND = 0 The solution is  $x = 0$ . No iterations were performed.
- IND = 1 The equations  $\bar{A}x = \bar{b}$  are probably compatible. A solution  $x$  has been obtained which is sufficiently accurate, given the values ATOL and BTOL.
- IND = 2 The equations  $\bar{A}x = \bar{b}$  are probably not compatible. A least squares solution  $x$  has been obtained which is sufficiently accurate, given the value ATOL.
- IND = 3 An estimate COND of  $\text{cond}(\bar{A})$  exceeds CONLIM. The vector  $x$  is the most recent approximation of a solution for  $\bar{A}x = \bar{b}$ .
- IND = 4 The equations  $\bar{A}x = \bar{b}$  are probably compatible. A solution  $x$  has been obtained which is as accurate as seems reasonable on this machine.
- IND = 5 The equations  $\bar{A}x = \bar{b}$  are probably not compatible. A least squares solution  $x$  has been obtained which is as accurate as seems reasonable on this machine.
- IND = 6  $\text{cond}(\bar{A})$  appears to be so large that there is not much point in doing further iterations. The vector  $x$  is the most recent approximation of a solution for  $\bar{A}x = \bar{b}$ .
- IND = 7 MXITER iterations were performed. More iterations are needed. The vector  $x$  is the most recent approximation of a solution for  $\bar{A}x = \bar{b}$ .

#### Remarks.

- (1) A large estimate of the condition number  $\text{cond}(\bar{A})$  may be due to rank deficiency or near rank deficiency of the matrix  $\bar{A}$ . If it is suspected that a large estimate of  $\text{cond}(\bar{A})$  has occurred for this reason, then it is recommended that CONLIM be set to a moderate value such as  $\sqrt{\epsilon}$  where  $\epsilon$  is the smallest value such that  $1 + \epsilon > 1$  ( $\epsilon = 2^{-47}$  for the CDC 6000-7000 series computers). Setting CONLIM to 0 is equivalent to setting CONLIM to  $\epsilon^{-1}$ .
- (2) The vector  $b$  is the only input argument modified by the routine.

**Algorithm.** BLSQ employs an iterative algorithm developed by Golub and Kahan.

**Programming.** BLSQ calls the subroutines NORMLZ, BVPRD1, BTPRD1, SCOPY, and SSCAL. The function SNRM2 is also used. BSLQ is an adaptation by A. H. Morris of the subroutine LSQR, written by Christopher C. Paige ( McGill University, Montreal, Canada) and Michael A. Saunders (Stanford University).

#### References.

- (1) Paige, C. C. and Saunders, M. A., "LSQR: An Algorithm for Sparse Linear Equations and Sparse Least Squares," *ACM Trans. Math Software* **8** (1982), pp. 43-71.
- (2) \_\_\_\_\_, "Algorithm 583. LSQR: Sparse Linear Equations and Least Squares Problems," *ACM Trans. Math Software* **8** (1982), pp. 195-209.

## ITERATIVE LEAST SQUARES SOLUTION OF SPARSE LINEAR EQUATIONS

Given an  $m \times n$  matrix  $A$ , a column vector  $b$  of dimension  $m$ , and a real number  $\lambda$ . Let  $\bar{A} = \begin{pmatrix} A \\ \lambda I \end{pmatrix}$  where  $I$  is the  $n \times n$  identity matrix, and let  $\bar{b} = \begin{pmatrix} b \\ 0 \end{pmatrix}$ . The problem is to find a column vector  $x$  of dimension  $n$  which is a least squares solution of  $\bar{A}x = \bar{b}$ . If  $A$  is sparse then the following subroutines are available for solving this problem.

```
CALL SPLSQ( $m, n, A, IA, JA, \lambda, b, x, ATOL, BTOL, CONLIM, MXITER,$   
            $IND, ITER, COND, RNORM, XNORM, WK$ )  
CALL STLSQ( $m, n, TA, ITA, JTA, \lambda, b, x, ATOL, BTOL, CONLIM, MXITER,$   
            $IND, ITER, COND, RNORM, XNORM, WK$ )
```

If SPLSQ is called then  $A, IA, JA$  are arrays containing the matrix  $A$  in sparse form. Otherwise, if STLSQ is called then  $TA, ITA, JTA$  are arrays containing the transpose matrix  $A^t$  in sparse form. An iterative procedure is used to obtain a least squares solution  $x$  of  $\bar{A}x = \bar{b}$ . The vector  $b$  is modified by the routines.

ATOL and BTOL are input arguments which specify the relative accuracy of  $A$  and  $b$  respectively. For example, if it is estimated that  $b$  is accurate to  $k$  decimal digits then one may set  $BTOL = 10^{-k}$ . It is required that  $ATOL \geq 0$  and  $BTOL \geq 0$ . If  $ATOL = 0$  or  $BTOL = 0$ , then it is assumed that  $A$  or  $b$  is accurate to machine precision.

Let  $\text{cond}(\bar{A})$  denote the condition number of  $\bar{A}$  relative to the Frobenius norm.<sup>1</sup> In each iteration of the algorithm being used, an estimate is made of the condition number  $\text{cond}(\bar{A})$ . The estimates form a monotonically nondecreasing sequence. The input argument CONLIM is an upper limit on  $\text{cond}(\bar{A})$ . If  $\text{CONLIM} > 0$  then the routines terminate when an estimate of  $\text{cond}(\bar{A})$  exceeds CONLIM. This termination may be needed to prevent small or zero singular values of  $\bar{A}$  from coming into effect and causing damage to the solution  $x$ . CONLIM may be ignored by being set to 0. It is assumed that  $\text{CONLIM} \geq 0$ .

The input argument MXITER is the maximum number of iterations that are permitted. Normally the routines require less than  $4n$  iterations. The related argument ITER is a variable. When the routines terminate  $ITER =$  the number of iterations that were performed.

COND, RNORM, and XNORM are variables. When the routines terminate  $\text{COND} =$  the last estimate made for  $\text{cond}(\bar{A})$ ,  $\text{RNORM} = \|\bar{A}x - \bar{b}\|$ , and  $\text{XNORM} = \|x\|$ .<sup>2</sup>

WK is an array of dimension  $2n$  or larger that is a work space for the routines.

<sup>1</sup> $\text{cond}(\bar{A}) = \|\bar{A}\|_F \|\bar{A}^+\|_F$  where  $\bar{A}^+$  is the pseudoinverse of  $\bar{A}$ . Here  $\|C\|_F = \sqrt{\sum c_{ij}^2}$  for any matrix  $C = (c_{ij})$ .

<sup>2</sup> $\|c\| = \sqrt{\sum c_i^2}$  for any vector  $c = (c_1, c_2, \dots)$ .

The equations  $\bar{A}x - \bar{b}$  are considered to be *compatible* if for any least squares solution  $x$ ,  $\|\bar{A}x - \bar{b}\| = 0$ . IND is a variable that reports the status of the results. When the routines terminate, IND has one of the following values:

- IND = 0 The solution is  $x = 0$ . No iterations were performed.
- IND = 1 The equations  $\bar{A}x - \bar{b}$  are probably compatible. A solution  $x$  has been obtained which is sufficiently accurate, given the values ATOL and BTOL.
- IND = 2 The equations  $\bar{A}x - \bar{b}$  are probably not compatible. A least squares solution  $x$  has been obtained which is sufficiently accurate, given the value ATOL.
- IND = 3 An estimate COND of  $\text{cond}(\bar{A})$  exceeds CONLIM. The vector  $x$  is the most recent approximation of a solution for  $\bar{A}x - \bar{b}$ .
- IND = 4 The equations  $\bar{A}x - \bar{b}$  are probably compatible. A solution  $x$  has been obtained which is as accurate as seems reasonable on this machine.
- IND = 5 The equations  $\bar{A}x - \bar{b}$  are probably not compatible. A least squares solution  $x$  has been obtained which is accurate as seems reasonable on this machine.
- IND = 6  $\text{cond}(\bar{A})$  appears to be so large that there is not much point in doing further iterations. The vector  $x$  is the most recent approximation of a solution for  $\bar{A}x - \bar{b}$ .
- IND = 7 MXITER iterations were performed. More iterations are needed. The vector  $x$  is the most recent approximation of a solution for  $\bar{A}x - \bar{b}$ .

#### Remarks.

- (1) A large estimate of the condition number  $\text{cond}(\bar{A})$  may be due to rank deficiency or near rank deficiency of the matrix  $\bar{A}$ . If it is suspected that a large estimate of  $\text{cond}(\bar{A})$  has occurred for this reason, then it is recommended that CONLIM be set to a moderate value such as  $\sqrt{\epsilon}$  where  $\epsilon$  is the smallest value such that  $1 + \epsilon > 1$  ( $\epsilon = 2^{-47}$  for the CDC 6000-7000 series computers). Setting CONLIM to 0 is equivalent to setting CONLIM to  $\epsilon^{-1}$ .
- (2) The vector  $b$  is the only input argument modified by the routine.

**Algorithm.** SPLSQ and STLSQ employ an iterative algorithm developed by Golub and Kahan.

**Programming.** SPLSQ and STLSQ call the subroutines NORMLZ, MVPRD1, MTPRD1, SCOPY, and SSCAL. The function SNRM2 is also used. SPLSQ and STLSQ are adaptations by A. H. Morris of the subroutine LSQR, written by Christopher C. Paige (McGill University, Montreal, Canada) and Michael A. Saunders (Stanford University).

#### References.

- (1) Paige, C. C. and Saunders, M. A., "LSQR: An Algorithm for Sparse Linear Equations and Sparse Least Squares," *ACM Trans. Math Software* 8 (1982), pp. 43-71.
- (2) \_\_\_\_\_, "Algorithm 583. LSQR: Sparse Linear Equations and Least Squares Problems," *ACM Trans. Math Software* 8 (1982), pp. 195-209.

## MINIMIZATION OF FUNCTIONS OF A SINGLE VARIABLE

Let  $F(x)$  be a continuous real-valued function defined for  $a \leq x \leq b$ . Then the following subroutine is available for finding a local minimum of  $F(x)$ .

**CALL FMIN**( $F, a, b, x, w, \text{AERR}, \text{RERR}, \text{ERROR}, \text{IND}$ )

It is assumed that  $a \leq b$ . FMIN finds a value  $x$  in the interval  $[a, b]$  which is a local minimum of  $F$ . ERROR and  $w$  are variables. When FMIN terminates,  $w = F(x)$  and ERROR is the estimated maximum absolute error of  $x$ .

The input arguments AERR and RERR are the absolute and relative error tolerances to be satisfied. For example, if  $k$  significant digit accuracy is desired then one may set  $\text{RERR} = 10^{-k}$ . It is assumed that  $\text{AERR} \geq 0$  and  $\text{RERR} \geq 0$ . The setting  $\text{AERR} = 0$  is equivalent to the setting  $\text{AERR} = 10^{-20}$ , and the setting  $\text{RERR} = 0$  is a request for machine precision.

IND is a variable that reports the status of the results.  $\text{IND} = 0$  if  $x$  is found to the desired accuracy. Otherwise,  $\text{IND} = 1$  when  $x$  cannot be obtained to the desired accuracy. In this case,  $w$  satisfies the tolerances AERR and RERR.

**Note.**  $F$  must be declared in the calling program to be of type EXTERNAL.

**Algorithm.** The golden section search procedure is used.

**Programming.** The function SPMPAR is called. FMIN was written by A. H. Morris.

## MINIMIZATION OF FUNCTIONS OF N VARIABLES

Let  $f(x)$  be a real-valued function of  $n$  variables  $x = (x_1, \dots, x_n)$  where  $n \geq 2$ . If  $f(x)$  is twice continuously differentiable then the following subroutine is available for finding a local minimum of  $f(x)$ .

**CALL OPTF( $F, n, RERR, ITER, X, FVAL, IND, WK$ )**

$X$  is an array of dimension  $n$  and  $FVAL$  a variable. On input,  $X$  contains an initial guess  $a = (a_1, \dots, a_n)$  to a minimum of  $f$ . When OPTF terminates,  $X$  contains the final estimate  $x = (x_1, \dots, x_n)$  of a local minimum of  $f$  and  $FVAL = f(x)$ .

The argument  $F$  is the name of a user defined subroutine that has the format:

**CALL  $F(n, X, FVAL)$**

Here  $X$  is an array of dimension  $n$  containing a point  $x = (x_1, \dots, x_n)$ , and  $FVAL$  is a variable.  $F$  sets  $FVAL$  to the value of the function  $f$  at the point  $x$ .  $F$  must be declared in the calling program to be of type EXTERNAL.

$RERR$  is an input argument that specifies the relative accuracy of  $F$ . If it is estimated that the subroutine  $F$  produces results accurate to  $k$  significant digits then one may set  $RERR = 10^{-k}$ . It is required that  $RERR \geq 0$ . If  $RERR = 0$  then it is assumed that  $F$  produces results accurate to machine precision.

When OPTF is called, line search iteration is performed to find the local minimum of  $f$ .  $ITER$  is a variable. On input,  $ITER$  is the maximum number of iterations that are permitted. When the routine terminates,  $ITER$  = the number of iterations that were actually performed.

$WK$  is an array of dimension  $n(n+8)$  or larger that is a work space for the routine.

$IND$  is a variable that reports the status of the results. When the routine terminates,  $IND$  has one of the following values:

- $IND = -1$  (Input error)  $n \leq 0$ .
- $IND = -2$  (Input error)  $n = 1$ .
- $IND = -4$  (Input error)  $ITER \leq 0$ .
- $IND = -5$  (Input error) Either  $RERR < 0$  or  $RERR > 10^{-4}$ .
- $IND = 1$  A local minimum  $x$  was found. The gradient of  $f$  at  $x$  was considered to be sufficiently small.
- $IND = 2$  The steps taken became so small that OPTF had to terminate.  $X$  is probably a local minimum, but it need not be a local minimum. The algorithm frequently requires exceedingly small steps to be taken, no matter whether  $X$  is close to or far from a local minimum.
- $IND = 3$  A local minimum has possibly been found. OPTF could not find a point for which  $f$  would take a smaller value.
- $IND = 4$   $ITER$  iterations were performed.
- $IND = 5$  The algorithm appears to be diverging. This generally occurs when  $f$  is unbounded from below.

When an input error is detected, the routine immediately terminates.

**Accuracy and Efficiency.** OPTF is frequently extremely efficient in finding a value FVAL which roughly approximates a local minimum value  $f(x_0)$ , but at times it can be quite slow in obtaining FVAL to greater precision. A rough approximation is often obtained in 20–30 iterations. If  $f(x_0) \neq 0$  then FVAL may be accurate to 4–6 significant digits after 20–30 additional iterations, or FVAL may not be accurate to 1 significant digit after several hundred iterations. 4–6 digit accuracy is the greatest precision that can be expected. In general, it is recommended that  $\text{ITER} \leq 200$ . Each iteration can take considerable time, even if the subroutine  $F$  is cheap to evaluate.

**Remark.** OPTF can be quite sensitive to the scaling of the variables  $(x_1, \dots, x_n)$ . The routine tends to operate more efficiently when the components of a local minimum  $x = (x_1, \dots, x_n)$  are all roughly of the same magnitude. If the components are of considerably different magnitudes (say  $|x_1| \approx 10^{-6}$  and  $|x_2| \approx 10^3$ ) then convergence may be extremely slow. In such a case, OPTF attempts to rescale the variables, but the rescaling is not always helpful.

**Algorithm.** The line search algorithm given in pp. 325–327 of the reference is employed. Also, BFGS secant updates for the hessian are used.

**Programming.** OPTF employs the subroutines OPTDRV, OPCHK1, OPSTP, FXDEC, SCALEX, LLTSLV, FSTOFD, FSTOCD, LNSRCH, SECFAC, QRUPDT, and JROT. OPTF, FXDEC, and SCALEX were written by A.H. Morris. The remaining subroutines were written by Robert B. Schnabel (University of Colorado at Boulder) and modified by A.H. Morris. The functions SDOT, SNRM2, and SPMPAR are also used.

**Reference.** Dennis, J.E. and Schnabel, R.B., *Numerical Methods for Unconstrained Optimization and Nonlinear Equations*, Prentice-Hall, Englewood Cliffs, NJ, 1983.

## UNCONSTRAINED MINIMUM OF THE SUM OF SQUARES OF NONLINEAR FUNCTIONS

Let  $f_1(x), \dots, f_m(x)$  be  $m$  real-valued functions of  $n$  real variables  $x = (x_1, \dots, x_n)$  where  $m \geq n$ . The problem under consideration is to find a point  $x$  which minimizes the function  $\phi(x) = \sum_{i=1}^m f_i(x)^2$ . Assume that each  $f_i(x)$  is differentiable and that an initial guess  $a = (a_1, \dots, a_n)$  to a minimum of  $\phi(x)$  is given. Then the following subroutine is available for finding a point which minimizes  $\phi(x)$ .

**CALL LMDIFF( $F, m, n, X, FVEC, EPS, TOL, INFO, IWK, WK, \ell$ )**

$X$  is an array of dimension  $n$  and  $FVEC$  an array of dimension  $m$ . On input  $X$  contains the starting point  $a = (a_1, \dots, a_n)$ . When LMDIFF terminates,  $X$  contains the final estimate  $x = (x_1, \dots, x_n)$  of a minimum of  $\phi$  and  $FVEC$  contains the values of the functions  $f_1, \dots, f_m$  at the output point in  $X$ .

The argument  $F$  is the name of a user defined subroutine that has the format:

**CALL  $F(m, n, X, FVEC, IFLAG)$**

Here  $X$  is an array of dimension  $n$ ,  $FVEC$  an array of dimension  $m$ , and  $IFLAG$  an integer variable. The array  $X$  contains a point  $x = (x_1, \dots, x_n)$ . Normally  $F$  evaluates the functions  $f_1, \dots, f_m$  at this point and stores the results in  $FVEC$ . However, if  $x$  does not lie in the domain of  $f_1, \dots, f_m$  then this cannot be done. In this case, the argument  $IFLAG$  (which will have been assigned a nonnegative value by LMDIFF) should be reset by  $F$  to a negative value. This will signal LMDIFF to terminate.  $F$  must be declared in the calling program to be of type EXTERNAL.

$EPS$  is an input argument which specifies the relative accuracy of  $F$ . If it is estimated that the subroutine  $F$  produces results accurate to  $k$  significant decimal digits then one may set  $EPS = 10^{-k}$ . It is required that  $EPS \geq 0$ . If  $EPS = 0$  then it is assumed that  $F$  produces results accurate to machine precision.

$TOL$  is an input argument which specifies the desired accuracy to be attained. The Euclidean norm  $\|x\| = \sqrt{\sum_i x_i^2}$  is employed. If  $\bar{x}$  denotes an actual minimum of  $\phi$ , then LMDIFF terminates when an iterate  $x$  is generated for which it is estimated that

$$(1) \phi(x) \leq (1 + TOL)^2 \phi(\bar{x}) \text{ or}$$

$$(2) \|D(x - \bar{x})\| \leq TOL \cdot \|D\bar{x}\|$$

is satisfied. In (2)  $x$  and  $\bar{x}$  are regarded as column vectors, and  $D$  is a diagonal matrix generated by LMDIFF whose entries are scaling factors. For convenience, criterion (1) is called the  $F$ -convergence (or  $\phi$ -convergence) test and criterion (2) is called the  $x$ -convergence test. It is required that  $TOL \geq 0$ . In order for the convergence tests to work properly, it is recommended that  $TOL$  always be smaller than  $10^{-5}$ .

$IWK$  is an array of dimension  $n$  and  $WK$  is an array of dimension  $\ell$ .  $IWK$  and  $WK$  are work spaces. It is assumed that  $\ell \geq mn + 5n + m$ .



INFO is an integer variable that reports the status of the results. When LMDIFF terminates, INFO has one of the following values:

- INFO < 0 This occurs when the user terminates the execution of LMDIFF by resetting the argument IFLAG in the subroutine *F* to a negative value. Then INFO = the negative value of IFLAG.
- INFO = 0 (Input error)  $1 \leq n \leq m$ ,  $\text{EPS} \geq 0$ ,  $\text{TOL} \geq 0$ , or  $\ell \geq mn + 5n + m$  is violated.
- INFO = 1 The *F*-convergence test has been satisfied.
- INFO = 2 The *x*-convergence test has been satisfied.
- INFO = 3 The *F*-convergence and *x*-convergence tests have been satisfied.
- INFO = 4 The gradient of  $\phi$  is 0 at point *X*.
- INFO = 5 The number of calls to the subroutine *F* has reached or exceeded  $200(n + 1)$ .
- INFO = 6 TOL is too small. No further reduction in the value of  $\phi(x)$  is possible.
- INFO = 7 TOL is too small. No further improvement in the accuracy of *X* is possible.

When LMDIFF terminates, if INFO  $\neq 0$  then *X* contains the final iterate that was generated. Also, if INFO  $\geq 1$  then FVEC contains the values of the functions  $f_1, \dots, f_m$  at this iterate. If INFO = 4 then *X* should be examined very closely. The gradient of  $\phi$  can be 0 when *X* is a local minimum or maximum, or when *X* is a saddle point. If INFO = 5 then it may (or may not) be helpful to continue the procedure by recalling LMDIFF with the current point in *X* as the new starting point. Since TOL is a relative tolerance, this setting can occur when  $\phi(0) = 0$ .

**Algorithm.** A modified form of the Levenberg-Marquardt algorithm is employed.

**Programming.** LMDIFF is a slightly modified version of the MINPACK-1 subroutine LMDIF1. The MINPACK-1 subroutines LMDIF, SPMPAR, ENORM, FDJAC2, LMPAR, QRFAC, and QRSOLV are employed. The subroutines were written by Jorge J. More, Burton S. Garbow, and Kenneth E. Hillstom (Argonne National Laboratory).

#### References.

- (1) More, J. J., Garbow, B. S., and Hillstom, K. E., *User Guide for MINPACK-1*, Argonne National Laboratory Report ANL-80-74, Argonne, Illinois, 1980.
- (2) More, J. J., "The Levenberg-Marquardt Algorithm: Implementation and Theory," *Numerical Analysis*, G. A. Watson (ed.), Springer-Verlag, 1977.

## LINEAR PROGRAMMING

Let  $A = (a_{ij})$  be an  $m \times n$  matrix,  $B$  an array containing  $b_1, \dots, b_m$ , and  $C$  an array containing  $c_1, \dots, c_n$  where  $a_{ij}, b_i, c_j$  are real. Consider the problem of finding nonnegative values  $x_1, \dots, x_n$  which maximize or minimize the function  $c_1x_1 + \dots + c_nx_n$  subject to the constraints:

$$\begin{aligned} a_{11}x_1 + \dots + a_{1n}x_n \{ \leq, =, \geq \} b_1 \\ \vdots \\ a_{m1}x_1 + \dots + a_{mn}x_n \{ \leq, =, \geq \} b_m \end{aligned}$$

In each constraint

$$a_{i1}x_1 + \dots + a_{in}x_n \{ \leq, =, \geq \} b_i$$

only one of the relations  $\leq, =, \geq$  is used, but the relation may vary from constraint to constraint. The following subroutines are available for solving this problem.

```
CALL SMPLX(A, B, C, ka, m, n, IND, IBASIS, X, z, ITER, MXITER,
           NUMLE, NUMGE, BI, WK, IWK)
CALL SSPLX(TA, ITA, JTA, B, C, m, n, IND, IBASIS, X, z, ITER, MXITER,
           NUMLE, NUMGE, BI, WK, IWK)
```

It is assumed that  $m \geq 2, n \geq 2$ , and that each  $b_i \geq 0$ . If SMPLX is called then  $ka$  is the number of rows in the dimension statement for  $A$  in the calling program. Otherwise, if SSPLX is called then  $TA, ITA, JTA$  are arrays containing the transpose matrix  $A^t$  in sparse form.

The constraints  $a_{i1}x_1 + \dots + a_{in}x_n \{ \leq, =, \geq \} b_i$  are assumed to be ordered so that the  $\leq$  constraints are followed by the  $\geq$  constraints, and the  $=$  constraints come last. NUMLE and NUMGE have the values:

NUMLE = the number of  $\leq$  constraints.

NUMGE = the number of  $\geq$  constraints.

It is assumed that  $NUMLE \geq 0, NUMGE \geq 0$ , and  $NUMLE + NUMGE \leq m$ .

When SMPLX or SSPLX is called, the routine attempts to maximize  $\sum_j c_j x_j$  subject to the constraints. A modified form of the primal simplex algorithm is employed. Frequently the procedure requires less than  $5m$  iterations to perform the task. The argument MXITER has the value:

MXITER = the maximum number of iterations that may be performed.

This argument is provided by the user. The related argument ITER is a variable that is set by the routine. When the routine terminates, ITER has for its value the number of iterations that were performed.

IND is a variable and IBASIS an array of dimension  $m$ . IBASIS contains the indices  $i$  of the current basic variables  $x_i$ , and IND is used for input/output purposes. On input IND is normally set by the user to 0. If  $IND = 0$  then the routine selects its own

beginning basis and stores the appropriate indices in IBASIS. [The remainder of this paragraph may be skipped by anyone not acquainted with the simplex algorithm.] If the user wishes to use his own beginning basis, then IND must be set to 1 and the indices of the initial basic variables stored in IBASIS. *It is not required that the initial basis be selected so that the basic variables are nonnegative.* The initial basic variables may be original, slack, surplus, or artificial variables. Slack and surplus variables are automatically provided for the  $\leq$  and  $\geq$  constraints, and artificial variables for the  $=$  constraints. The routine defines  $x_{n+i}$  to be the slack, surplus, or artificial variable for the  $i^{\text{th}}$  constraint  $a_{i1}x_1 + \dots + a_{in}x_n \{ \leq, =, \geq \} b_i$ . If  $\text{IND} = 0$  then the slack, surplus, and artificial variables are the initial basic variables that are employed.

On output IND reports the status of the results. The routine assigns IND one of the following values:

- IND = 0 The problem was solved.
- IND = 1 The problem has no solution.
- IND = 2 MXITER iterations were performed. More iterations are needed.
- IND = 3 Sufficient accuracy cannot be maintained to solve the problem.
- IND = 4 The problem has an unbounded solution.
- IND = 5 An input error was detected. (See below)
- IND = 6 A possible solution was obtained. The routine is not certain if the solution is correct.

$X$  is an array of dimension  $n + \text{NUMLE} + \text{NUMGE}$  and  $z$  is a variable. If  $\text{IND} = 0$  or  $\text{IND} = 6$ , then  $z$  has for its value the maximum value obtained for  $\sum_j c_j x_j$  and  $X$  contains the values obtained for the original, slack, and surplus variables. If  $\text{IND} \neq 5$  then IBASIS contains the indices of the basic variables currently in effect when the routine terminates.

BI is an array of dimension  $m^2$  that is used for storing the inverse of the basis matrix. The order of the column vectors of the basis matrix corresponds to the order of the basic variables given in IBASIS. If  $\text{IND} \neq 3, 5$  on output then BI contains the inverse of the basis matrix currently in effect when the routine terminates.

WK is an array of dimension  $2m$  or larger, and IWK is an array of dimension  $2m + n$  or larger. WK and IWK are work spaces.

**Input Errors.**  $\text{IND} = 5$  occurs on output when one of the following conditions is violated:

- (1)  $n \geq 2$  and  $ka \geq m \geq 2$ .
- (2)  $\text{NUMLE} + \text{NUMGE} \leq m$ .
- (3) Each  $b_i \geq 0$ .
- (4) The basis matrix specified by the user in IBASIS (when  $\text{IND} = 1$  on input) is nonsingular and sufficiently well conditioned so that its inverse can be computed.

#### Remarks.

- (1)  $A, B, C, \text{TA}, \text{ITA}, \text{JTA}$  are not modified by the routines.
- (2) The routines maximize  $\sum_j c_j x_j$ . This function can be minimized by maximizing  $\sum_j (-c_j) x_j$  and then changing the sign of the result.

- (3) SMPLX and SSPLX generate the same results. For efficiency, SSPLX should be used when  $A$  is sparse.

**Algorithm.** A three step procedure is used. The first step eliminates the negative variables. Then phases (1) and (2) of the primal simplex algorithm are invoked. Negative variables are eliminated as follows: Let  $x_{B1}, \dots, x_{Bm}$  be the basic variables and  $y_{ij}$  the components of the simplex tableau.

- (1) Compute  $d_j = \sum_i y_{ij}$  for each nonbasic variable  $x_j$  where the sum  $\sum_i$  is for all  $i$  where  $x_{Bi} < 0$ . If all  $d_j \geq 0$  then the problem has no feasible solution. Otherwise, select  $k$  so that  $d_k = \min_j d_j$ . Then  $x_k$  is the variable to be made basic.
- (2) If  $x_{Bj} \geq 0$  and  $y_{jk} > 0$  for some  $j$  then go to (3). Otherwise, select a negative variable  $x_{Br}$  to become nonbasic where  $x_{Br}/y_{rk} = \max\{x_{Bj}/y_{jk} : x_{Bj} < 0 \text{ and } y_{jk} < 0\}$ . Then update the basis and go to (5).
- (3) Compute  $\epsilon = \min\{x_{Bj}/y_{jk} : x_{Bj} \geq 0 \text{ and } y_{jk} > 0\}$  and check if a negative variable  $x_{Bj}$  exists that satisfies the conditions:  
 (\*)  $y_{jk} < 0$  and  $\epsilon \geq x_{Bj}/y_{jk}$   
 If such a variable exists then go to (4). Otherwise, select a nonnegative variable  $x_{Br}$  to become nonbasic where  $y_{rk} > 0$  and  $x_{Br}/y_{rk} = \epsilon$ . Then update the basis and go to (1).
- (4) Select a negative variable  $x_{Br}$  to become nonbasic where  $x_{Br}/y_{rk} = \max\{x_{Bj}/y_{jk} : x_{Bj} < 0 \text{ and } x_{Bj} \text{ satisfies } (*)\}$ . Then update the basis and go to (5).
- (5) Check if there are any remaining negative variables. If not, then we are finished. Otherwise go to (1).

**Programming.** SMPLX and SSPLX employ the subroutines SMPLX1, SSPLX1, and CROUT1. These routines were written by A. H. Morris. The function SPMPAR is also used.

## THE ASSIGNMENT PROBLEM

Let  $C = (c_{ij})$  be an  $n \times n$  matrix (the cost matrix). The problem under consideration is to find an  $n \times n$  matrix  $x = (x_{ij})$  which minimizes  $T = \sum_{i=1}^n \sum_{j=1}^n c_{ij}x_{ij}$  and satisfies:

$$(1) \sum_{i=1}^n x_{ij} = 1 \text{ for } j = 1, \dots, n$$

$$(2) \sum_{j=1}^n x_{ij} = 1 \text{ for } i = 1, \dots, n$$

$$(3) \text{ Each } x_{ij} = 0 \text{ or } 1$$

Each  $x$  which satisfies (1)–(3) is called an *assignment*. For each such  $x$ , from (1) and (3) we note that for each  $j$  there exists a unique integer  $\pi(j)$  such that  $x_{\pi(j),j} = 1$ . Also, (2) and (3) assert that  $\pi$  is a permutation of  $\{1, \dots, n\}$ . Conversely, for any permutation  $\pi$  there corresponds an assignment  $x$  defined by  $x_{\pi(j),j} = 1$  and  $x_{ij} = 0$  for  $i \neq \pi(j)$ . Thus, the problem is to find a permutation  $\pi$  of  $\{1, \dots, n\}$  which minimizes  $T = \sum_{j=1}^n c_{\pi(j),j}$ . The following subroutine is available for solving this problem when all  $c_{ij}$  are integers.

**CALL ASSGN( $n, C, JC, T, IWK, IERR$ )**

$C$  is a 2-dimensional integer array of dimension  $n \times (n + 1)$ ,  $JC$  an integer array of dimension  $n$ , and  $T$  an integer variable. It is assumed that  $n \geq 2$  and that the first  $n$  columns of  $C$  contain the cost matrix  $(c_{ij})$ . [The  $(n + 1)^{st}$  column of  $C$  is a work space for the routine.] When ASSGN is called, the desired permutation  $\pi$  is obtained and the values  $\pi(1), \dots, \pi(n)$  stored in  $JC$ . Also  $T$  is assigned the minimized value  $\sum_j c_{\pi(j),j}$ .

$IWK$  is an array of dimension  $7n + 2$  or larger that is a work space for the routine.

$IERR$  is a variable that is set by the routine. If  $JC$  and  $T$  are obtained then  $IERR$  is assigned the value 0. Otherwise, if the problem cannot be solved because of integer overflow, then  $IERR = 1$ .

### Remarks.

- (1)  $C$  is destroyed by the routine.
- (2) ASSGN minimizes  $T = \sum_j c_{\pi(j),j}$ . This function can be maximized by minimizing  $\sum_j (-c_{\pi(j),j})$  and then changing the sign of the result.

**Programming.** ASSGN calls the subroutine ASSGN1. ASSGN1 was written by Giorgio Carpaneto and Paolo Toth (University of Bologna, Italy), and modified by A. H. Morris. The function IPMPAR is also used.

**Reference.** Carpaneto, G., and Toth, P., "Algorithm 548, Solution of the Assignment Problem," *ACM Trans. Math Software* 6 (1980), pp. 104–111.

## 0-1 KNAPSACK PROBLEM

Given  $n \geq 2$  items, each having a profit  $p_j > 0$  and weight  $w_j > 0$ , and  $m \geq 1$  containers (knapsacks), each having a capacity  $k_i > 0$ . Let  $x_{ij} = 1$  if item  $j$  is assigned to knapsack  $i$ . Otherwise, let  $x_{ij} = 0$ . Then the problem is to find an assignment  $x_{ij}$  of the items to the knapsacks which maximizes  $\sigma = \sum_{i=1}^m \sum_{j=1}^n p_j x_{ij}$  subject to

$$(1) \sum_{j=1}^n w_j x_{ij} \leq k_i \text{ for } i = 1, \dots, m, \text{ and}$$

$$(2) \sum_{i=1}^m x_{ij} \leq 1 \text{ for } j = 1, \dots, n.$$

Condition (2) states that each item may be assigned to a (single) knapsack or be rejected. It can be assumed, without loss of generality, that

- (a) the knapsacks are ordered so that  $k_1 \leq \dots \leq k_m$ ,
- (b)  $\min\{w_1, \dots, w_n\} \leq k_1$ ,
- (c)  $\max\{w_1, \dots, w_n\} \leq k_m$ , and
- (d)  $\sum_{j=1}^n w_j > k_n$ .

Then the following subroutine is available for solving this problem when all  $p_j$ ,  $w_j$ , and  $k_i$  are positive integers.

**CALL MKP**( $n, m, P, W, K, \text{NBCK}, L, \sigma, \text{TEMP}, \text{ITEMP}, \text{NUM}$ )

$P$  and  $W$  are integer arrays containing  $p_1, \dots, p_n$  and  $w_1, \dots, w_n$ , and  $K$  an integer array containing  $k_1, \dots, k_m$ .  $L$  is an integer array of dimension  $n$  or larger, and  $\sigma$  an integer variable. When MKP is called, if no input errors are detected then the maximum value for  $\sigma$  is obtained. Also,  $L(j) = i$  if item  $j$  has been assigned to knapsack  $i$  ( $j = 1, \dots, n$ ), and  $L(j) = 0$  if item  $j$  does not appear in the solution (i.e., if  $x_{ij} = \dots = x_{mj} = 0$ ).

A depth-first tree search employing backtracking is used. If no bound is placed on the number of back tracks that may be performed, then the exact maximum  $\sigma$  is assured. However, if the number of back tracks must be restricted, then only an approximation to the maximum may be obtained. The argument NBCK is available for limiting the backtracking. NBCK is a variable. On input, if NBCK = -1 then no restrictions are placed on the backtracking. Otherwise, if NBCK  $\neq$  -1, then it is assumed that NBCK is the maximum number of back tracks that are permitted. When the routine terminates, NBCK = the number of back tracks that were actually performed.

TEMP is a real array of dimension  $n$  or larger, and ITEMP an integer array of dimension NUM. It is assumed that  $\text{NUM} \geq 5m + 14n + 4mn + 3$ . TEMP and ITEMP are work spaces for the routine. It should be noted that TEMP is the only argument of MKP that is not of integer type.

**Error return.** If an input error is detected then  $\sigma$  is set to one of the following values:

$$\sigma = -1 \text{ if } m < 1 \text{ or } n < 2.$$

- $\sigma = -2$  if some  $p_j, w_j$ , or  $k_i$  is not positive.
- $\sigma = -3$  if  $\min\{w_1, \dots, w_n\} > k_1$ ; i.e., if a knapsack cannot contain any item.
- $\sigma = -4$  if  $\max\{w_1, \dots, w_n\} > k_m$ ; i.e., if an item cannot fit in any knapsack.
- $\sigma = -5$  if  $\sum_j w_j \leq k_m$ ; i.e., if knapsack  $m$  can contain all the items.
- $\sigma = -7$  if the knapsacks are not ordered so that  $k_1 \leq \dots \leq k_m$ .
- $\sigma = -8$  if  $\text{NUM} < 5m + 14n + 4mn + 3$ .

**Backtracking.** For  $\text{NBCK} = -1$ , the time required for finding the exact maximum depends primarily on the value of  $m$ , and can increase quite dramatically for very small increases in the value of  $m$ . It is recommended that this setting never be used when  $m > 10$ . Instead, if  $m > 10$  then a setting of  $\text{NBCK} \leq 50$  frequently suffices.

**Programming.** MKP employs the subroutines MKP1, SIGMA1, PI1, PARC, SKNP, and SKNP1. The interface subroutine MKP was written by A.H. Morris. The remaining subroutines were written by Silvano Martello (University of Bologna) and Paolo Toth (University of Florence) and modified by A.H. Morris. The subroutine RISORT is also used.

**Reference.** Martello, S. and Toth, P., "Algorithm 632. A Program for the 0-1 Multiple Knapsack Problem," *ACM Trans. Math Software* 11 (1985), pp. 135-140.

## INVERSION OF THE LAPLACE TRANSFORM

Let  $f(t)$  be a complex-valued function that is continuous for  $t > 0$  except possibly at a countable set of values  $t_k$  having no finite limit points. Then for complex  $z$  the Laplace transform  $F(z)$  of  $f(t)$  is defined by

$$F(z) = \int_0^{\infty} e^{-zt} f(t) dt$$

when the integral converges. If the integral converges for  $\text{Re}(z) > c$  but does not exist for  $\text{Re}(z) < c$ , then  $c$  is called the *abscissa of convergence* of  $F(z)$ . If  $\int_0^{\infty} |f(t)|e^{-at} dt < \infty$  for some real constant  $a$ , then  $F(z)$  is analytic for all  $z$  where  $\text{Re}(z) > a$ . Also, for any point  $t$  for which  $f(t)$  is continuous and sufficiently well-behaved, the value  $f(t)$  can be obtained from  $F$  by the inversion formula

$$f(t) = \frac{1}{2\pi i} \lim_{T \rightarrow \infty} \int_{\alpha-iT}^{\alpha+iT} e^{\lambda t} F(\lambda) d\lambda \quad (i = \sqrt{-1})$$

for any  $\alpha > a$ . If  $f(t)$  is real for  $t > 0$ , then  $\overline{F(z)} = F(\bar{z})$ . Given a transform  $F(z)$  where  $\overline{F(z)} = F(\bar{z})$ , then the following subroutine is available for computing  $f(t)$  for  $t > 0$ .

**CALL LAINV(MO,FUN,t,AERR,RERR,Y,c,ERROR,NUM,IERR)**

It is assumed that  $\overline{F(z)} = F(\bar{z})$ . The argument FUN is the name of a user defined subroutine for computing  $F(z)$ . FUN has the format:

**CALL FUN(x,y,A,B)**

$A$  and  $B$  are variables. For real arguments  $x$  and  $y$ ,  $A$  and  $B$  are assigned the values  $A = \text{Re}[F(x + iy)]$  and  $B = \text{Im}[F(x + iy)]$ . FUN must be declared in the calling program to be of type EXTERNAL.

IERR is a variable that is used both for input and output purposes. If  $\text{IERR} \geq 0$  on input then it is assumed that the abscissa of convergence is known and that  $c =$  the abscissa of convergence. Otherwise, if  $\text{IERR} < 0$  then the abscissa of convergence must be computed. In this case,  $c$  is a variable. LAINV sets  $c$  to the value that is obtained for the abscissa of convergence.

MO is an integer which specifies the search procedure to be used for finding the abscissa of convergence  $c$  when  $\text{IERR} < 0$  on input. A two-pass procedure is employed when  $\text{MO} = 0$ , and a one-pass procedure when  $\text{MO} \neq 0$ . When all the singularities of  $F(z)$  are real, the two-pass procedure ( $\text{MO} = 0$ ) is almost always considerably more efficient. Conversely, if none of the singularities of  $F(z)$  are real, or if  $F(z)$  has complex singularities that are to the right of real singularities, then the one-pass procedure ( $\text{MO} \neq 0$ ) is more efficient.

It is assumed that  $t > 0$  and that  $Y$  is a variable. When LAINV is called  $Y$  is set to the value obtained for  $f(t)$ .

AERR and RERR are the absolute and relative error tolerances to be used in computing  $F(t)$  ( $\text{AERR} \geq 0$  and  $\text{RERR} \geq 0$ ). If one wants accuracy to  $k$  significant digits then set



$RERR = 10^{-4}$ . If  $RERR = 0$  then it is assumed that  $f(t)$  is to be computed to machine accuracy. LAINV attempts to find a value  $Y$  which satisfies  $|Y - f| < \max\{AERR, RERR \cdot f(t)\}$ .

ERROR and NUM are variables that are set by the routine. When LAINV terminates, ERROR is a rough estimate of the absolute error  $|Y - f(t)|$  and NUM is the number of calls that were made to the subroutine FUN.

When LAINV terminates, IERR reports the status of the results. IERR is assigned one of the following values:

- IERR = 0  $Y$  was obtained to the desired accuracy.
- IERR = 1  $Y$  was obtained, but it may not be accurate because of inaccuracy in the computation of  $c$ . This setting occurs only when IERR < 0 on input.
- IERR = 2  $Y$  could not be obtained, possibly because too much accuracy was requested. Increase AERR and RERR, and rerun the problem.
- IERR = 3  $Y$  could not be obtained, possibly because of inaccuracy in the computation of  $c$  or too much accuracy was requested. Increase AERR and RERR, and rerun the problem. This setting occurs only when IERR < 0 on input.
- IERR = 4 (Input error) The argument  $t$  is not positive.  $Y$  and ERROR are assigned the values 0 and 1.
- IERR = 5 The abscissa of convergence  $c$  could not be found in the interval  $[-10^4, 10^4]$ .  $Y$ ,  $c$ , and ERROR are assigned the values 0, 0, and 1. This setting occurs only when IERR < 0 on input.
- IERR = 6 The argument  $t$  is too large for  $f(t)$  to be computed.  $Y$  and ERROR are assigned the values 0 and 1.

#### Remarks.

- (1) Accuracy decreases when  $t$  is near a discontinuity of  $f(t)$ .
- (2) The calculation may lose accuracy or fail when  $F(t)$  is oscillatory.

**Algorithm.** Given  $c$ ,  $f(t)$  is computed by a modification of the subroutine DLAINV developed by R. Piessens and R. Huysmans, where the real Wynn  $\epsilon$ -algorithm has been replaced with the complex Wynn  $\epsilon$ -algorithm.

When IERR < 0,  $c$  is calculated by the subroutine ABCON or the subroutine ABCON1. In ABCON, which is a two-pass search procedure, the abscissa  $x_1$  of the rightmost singularity in the strip  $-10^4 \leq x \leq 10^4$ ,  $|y| \leq .01$  is first determined. Then the abscissa of the rightmost singularity in the half-plane  $\text{Re}(z) \geq x_1$  is found. In ABCON1 these calculations are combined into a single-pass procedure.

In ABCON and ABCON1, the function  $F(z)/(z + |x_1| + 1)$  is integrated along paths  $C_1$  and  $C_2$  defined as follows:  $C_1$  is the straight line segment from  $(x_1, 0)$  to  $(x_1, .01)$ , followed by the straight line segment from  $(x_1, .01)$  to  $(\infty, .01)$ , and  $C_2$  is the straight line segment from  $(x_1, \infty)$  to  $(x_1, 0)$ , followed by the straight line segment from  $(x_1, 0)$  to  $(\infty, 0)$ . The integral along  $C_1$  vanishes if no singularity lies to the right of  $x_1$  in the strip  $|y| \leq .01$ , and the integral along  $C_2$  vanishes if no singularity lies in  $\text{Re}(z) \geq x_1$ . Otherwise, these

integrals are nonzero in most applications. Simpson's rule suffices for integrating along the finite line segment from  $(x_1, 0)$  to  $(x_1, .01)$ .

**Example.** Let  $F(z) = 1/(1 + z^2)$ , in which case  $f(t) = \sin t$ . The following code may be used for computing  $f(t)$  at  $t = 1, 1.1, 1.2, \dots, 1.9$  and storing the results in the array  $W$ .

```

      REAL W(10)
      EXTERNAL FT
C
      AERR = 1.E-30
      RERR = 1.E-12
      IERR = -1
      T = 1.0
      DO 10 I = 1,10
          CALL LAINV (1, FT, T, AERR, RERR, W(I), C, ERR, N, IERR)
          IF (IERR .GT. 1) STOP
          T = T + 0.1
10  CONTINUE

```

Here FT may be defined by:

```

      SUBROUTINE FT(X, Y, A, B)
      COMPLEX Z,W
C
      Z = CMPLX(X,Y)
      W = 1.0/(1.0 + Z**2)
      A = REAL(W)
      B = AIMAG(W)
      RETURN
      END

```

**Programming.** LAINV employs the subroutines ABCON, ABCON1, SRCH, ACOND, XCOND, LAINV1, CQEXT, QAGI1, QAGIE1, QELG, QK15I1, QPSRT, CDIVID and functions ACONDF, ACONDG, XCONDX, XCONDY, SPMPAR, EXPARG. LAINV and ABCON were written by Andrew H. van Tuyl (NSWC) and modified by A.H. Morris. ABCON1 was written by A.H. Morris. LAINV1 was written by Robert Piessens and Rudi Huysmans (University of Leuven, Herverlee, Belgium) and modified by Andrew van Tuyl. QAGI1 is a modification of QAGI by Andrew van Tuyl and A.H. Morris.

**Reference.** Piessens, R. and Huysmans, R., "Algorithm 619. Automatic Numerical Inversion of the Laplace Transform," *ACM Trans. Math Software* 10 (1984), pp. 348-353.

## FAST FOURIER TRANSFORM

Let  $n$  be a positive integer and  $\theta_j = 2\pi j/n$  for  $j = 0, 1, \dots, n-1$ . For any complex valued functions  $f$  and  $g$  defined on the points  $\theta_j$  let  $(f, g) = \sum_{j=0}^{n-1} f(\theta_j) \overline{g(\theta_j)}$ . Then  $(f, g)$  is an inner product when  $f$  and  $g$  are regarded as functions defined only on  $\theta_j$ . Also  $e^{ij\theta}$  ( $j = 0, 1, \dots, n-1$ ) form an orthogonal set of functions where each  $e^{ij\theta}$  has norm  $\sqrt{n}$ .<sup>1</sup> Thus, if  $f$  is a function that is approximated by  $f(\theta) = \sum_{j=0}^{n-1} c_j e^{ij\theta}$  then each  $c_j = \frac{1}{n}(f(\theta), e^{ij\theta})$ . The mapping  $f(\theta_j) \mapsto c_j$  given by

$$c_j = \frac{1}{n} \sum_{k=0}^{n-1} f(\theta_k) e^{-2\pi i j k/n}$$

is called the *discrete Fourier transform* and its inverse

$$f(\theta_j) = \sum_{k=0}^{n-1} c_k e^{2\pi i j k/n}$$

the *inverse discrete Fourier transform*. The following subroutines are available for computing these transforms.

**CALL FFT(C, n,  $\ell$ , IERR)**  
**CALL FFT1(A, B, n,  $\ell$ , IERR)**

Let  $c_j = a_j + ib_j$  ( $j = 0, 1, \dots, n-1$ ) be the data to be transformed. If FFT is called then  $C$  is a complex array containing  $c_0, c_1, \dots, c_{n-1}$  (where  $C(j+1) = c_j$  for  $j < n$ ). Otherwise, if FFT1 is called then  $A$  and  $B$  are real arrays containing  $a_0, a_1, \dots, a_{n-1}$  and  $b_0, b_1, \dots, b_{n-1}$  respectively.

The argument  $\ell$  may have the values 1 or -1, and IERR is a variable. When FFT or FFT1 is called, if there are no input errors then IERR is set to 0 and

$$\hat{c}_j = \sum_{k=0}^{n-1} c_k e^{2\pi i \ell j k/n}$$

is computed. The results  $\hat{c}_j = \hat{a}_j + i\hat{b}_j$  replace the original data  $c_j = a_j + ib_j$  in  $C$  (or  $A$  and  $B$ ).

**Restrictions on the argument  $n$ .** When FFT and FFT1 are called,  $n$  is factored by the routine into its prime factors. It is assumed that the largest prime factor of  $n$  is  $\leq 23$ . If  $n = \mu^2 \bar{n}$  where  $\bar{n}$  is the square free portion of  $n$ , then it is further assumed that  $\bar{n} \leq 210$  whenever  $\bar{n}$  is a product of two or more primes.

<sup>1</sup>Throughout this section  $i = \sqrt{-1}$ .

**Error Return.** If an input error is detected then IERR is set as follows:

IERR = 1 if  $n < 1$ .

IERR = 2 if  $n$  has too many factors.

IERR = 3 if  $n$  has a prime factor greater than 23 or the square free portion of  $n$  is greater than 210.

IERR = 4 if  $\ell \neq \pm 1$ .

The setting IERR = 2 can occur only when  $n > 4251528$ .

**Remark.** The complex array  $C$  is interpreted by FFT as a real array of dimension  $2n$ . If this association is not permitted by the FORTRAN being employed then use FFT1.

**Programming.** FFT and FFT1 are interface routines for the subroutine SFFT, which was written by Richard C. Singleton (Stanford Research Institute).

**Reference.** Singleton, R. C., "An Algorithm for Computing the Mixed Radix Fast Fourier Transform," *IEEE Trans. Audio and Electroacoustics*, vol. AU-17 (1969), pp. 93-103.

## MULTIVARIATE FAST FOURIER TRANSFORM

Let  $n_1, \dots, n_m$  be positive integers. For any  $J = (j_1, \dots, j_m)$  where  $j_\nu = 0, \dots, n_\nu - 1$  ( $\nu = 1, \dots, m$ ) let  $\theta_J$  denote the point  $(2\pi j_1/n_1, \dots, 2\pi j_m/n_m)$ . Also, for any complex valued functions  $f$  and  $g$  defined on the points  $\theta_J$  let  $(f, g) = \sum_J f(\theta_J) \overline{g(\theta_J)}$ . Then  $(f, g)$  is an inner product when  $f$  and  $g$  are regarded as functions defined only on  $\theta_J$ . Also the functions  $\phi_J(\theta) = \exp(ij_1\theta^1) \cdots \exp(ij_m\theta^m)$  form an orthogonal set where each  $\phi_J$  has the norm  $\sqrt{n_1 \cdots n_m}$ .<sup>1</sup> Thus, if  $f$  is a function that is approximated by  $f = \sum_J c_J \phi_J$  then each  $c_J = \frac{1}{n_1 \cdots n_m} (f, \phi_J)$ . The mapping  $f(\phi_J) \mapsto c_J$  given by

$$c_J = \frac{1}{n_1 \cdots n_m} \sum_K f(\theta_K) \exp(-2\pi i j_1 k_1/n_1) \cdots \exp(-2\pi i j_m k_m/n_m)$$

is called the *discrete multivariate Fourier transform* and its inverse

$$f(\theta_J) = \sum_K c_K \exp(2\pi i j_1 k_1/n_1) \cdots \exp(2\pi i j_m k_m/n_m)$$

the *inverse discrete multivariate Fourier transform*. The sums  $\sum_K$  are for all  $K = (k_1, \dots, k_m)$  where  $k_\nu = 0, 1, \dots, n_\nu - 1$  ( $\nu = 1, \dots, m$ ). The following subroutines are available for computing these transforms.

CALL MFFT( $C, N, m, \ell, \text{IERR}$ )  
CALL MFFT1( $A, B, N, m, \ell, \text{IERR}$ )

Let  $c_J = a_J + ib_J$  be the data to be transformed where  $J = (j_1, \dots, j_m)$  for  $j_\nu = 0, 1, \dots, n_\nu - 1$  ( $\nu = 1, \dots, m$ ). If MFFT is called then  $C$  is a 1-dimensional complex array containing the values  $c_J$  where  $c_J = C(1 + j_1 + j_2 n_1 + j_3 n_1 n_2 + \cdots + j_m n_1 \cdots n_{m-1})$ . Otherwise, if MFFT1 is called then  $A$  and  $B$  are 1-dimensional real arrays containing the data  $a_J$  and  $b_J$  respectively.

**Note.** If MFFT is used and  $m = 2$  or  $3$ , then instead of having to store the  $m$ -dimensional data  $c_J$  into a 1-dimensional array  $C$ , the data may be stored in  $C$  where  $C$  is defined to be an  $m$ -dimensional array. If  $m = 2$  then  $C$  may be declared to be of dimension  $n_1 \times n_2$ , in which case  $C(j_1 + 1, j_2 + 1) = c_J$  for all  $J = (j_1, j_2)$ . Similarly, if  $m = 3$  then  $C$  may be declared to be of dimension  $n_1 \times n_2 \times n_3$ , in which case  $C(j_1 + 1, j_2 + 1, j_3 + 1) = c_J$  for all  $J = (j_1, j_2, j_3)$ . Similar comments hold for  $A$  and  $B$  if MFFT1 is employed.

$N$  is an array containing the integers  $n_1, \dots, n_m$ . The argument  $\ell$  may have the values 1 and  $-1$ , and IERR is a variable. When MFFT or MFFT1 is called, if there are no input errors then IERR is set to 0 and the transform

$$\hat{c}_J = \sum_K c_K \exp(2\pi i \ell j_1 k_1/n_1) \cdots \exp(2\pi i \ell j_m k_m/n_m)$$

is computed. The results  $\hat{c}_J = \hat{a}_J + i\hat{b}_J$  replace the original data  $c_J = a_J + ib_J$  in  $C$  (or  $A$  and  $B$ ).

<sup>1</sup>Throughout this section  $i = \sqrt{-1}$  and  $\theta = (\theta^1, \dots, \theta^m)$  denotes an arbitrary point.

**Restrictions on the arguments  $n_1, \dots, n_m$ .** When MFFT and MFFT1 are called, each  $n_\nu$  is factored by the routine into its prime factors. It is assumed that the largest prime factor of  $n_\nu$  is  $\leq 23$ . If  $n_\nu = \mu_\nu^2 \bar{n}_\nu$ , where  $\bar{n}_\nu$  is the square free portion of  $n_\nu$ , then it is further assumed that  $\bar{n}_\nu \leq 210$  whenever  $\bar{n}_\nu$  is a product of two or more primes.

**Error Return.** If an input error is detected then IERR is set as follows:

IERR = 1 if some  $n_\nu < 1$ .

IERR = 2 if some  $n_\nu$  has too many factors.

IERR = 3 if some  $n_\nu$  has a prime factor greater than 23 or the square free portion of some  $n_\nu$  is greater than 210.

IERR = 4 if  $\ell \neq \pm 1$ .

IERR = 5 if  $m \leq 0$ .

The setting IERR = 2 can occur only when some  $n_\nu > 4251528$ .

**Remark.** The complex array  $C$  of dimension  $n_1 \cdots n_m$  is interpreted by MFFT as a real array of dimension  $2n_1 \cdots n_m$ . If this association is not permitted by the FORTRAN being employed then use MFFT1.

**Programming.** MFFT and MFFT1 are interface routines for the subroutine SFFT, which was written by Richard C. Singleton (Stanford Research Institute).

**Reference.** Singleton, R. C., "An Algorithm for Computing the Mixed Radix Fast Fourier Transform," *IEEE Trans. Audio and Electroacoustics*, vol. AU-17 (1969), pp 93-103.

## DISCRETE COSINE AND SINE TRANSFORMS

Let  $n$  be a positive integer and  $\theta_\nu = (\nu + 1/2)\pi/n$  for  $\nu = 0, 1, \dots, n-1$ . For any real valued functions  $f$  and  $g$  defined on the points  $\theta_\nu$  let  $(f, g) = \sum_{\nu=0}^{n-1} f(\theta_\nu) g(\theta_\nu)$ . Then  $(f, g)$  is an inner product when  $f$  and  $g$  are regarded as functions defined only on  $\theta_\nu$ . Also  $\cos j\theta$  ( $j = 0, 1, \dots, n-1$ ) form an orthogonal set of functions where  $\cos j\theta$  has norm  $\sqrt{n}$  when  $j = 0$  and norm  $\sqrt{n/2}$  when  $j \geq 1$ . Thus, if  $f$  is a function that is approximated by  $f(\theta) = a_0 + 2 \sum_{j=1}^{n-1} a_j \cos j\theta$  then each  $a_j = \frac{1}{n}(f(\theta), \cos j\theta)$ . The mapping  $f(\theta_\nu) \mapsto a_j$  is called the *discrete cosine transform* and its inverse  $a_j \mapsto f(\theta_\nu)$  the *inverse discrete cosine transform*.

Alternatively, the functions  $\sin j\theta$  for  $j = 1, \dots, n$  also form an orthogonal set where  $\sin j\theta$  has norm  $\sqrt{n/2}$  when  $j < n$  and norm  $\sqrt{n}$  when  $j = n$ . Thus, if  $f$  is a function that is approximated by  $f(\theta) = 2 \sum_{j=1}^{n-1} b_j \sin j\theta + b_n \sin n\theta$  then each  $b_j = \frac{1}{n}(f(\theta), \sin j\theta)$ . The mapping  $f(\theta_\nu) \mapsto b_j$  is called the *discrete sine transform* and its inverse  $b_j \mapsto f(\theta_\nu)$  the *inverse discrete sine transform*.

The subroutines COSQB and COSQF are available for computing the discrete cosine transform and its inverse, and the subroutines SINQB and SINQF are available for computing the discrete sine transform and its inverse. The subroutine COSQI provides information that is needed for the cosine and sine transform routines.

### CALL COSQI( $n$ , WK)

WK is an array of dimension  $3n + 15$  or larger that is a work space for the routines COSQB, COSQF, SINQB, and SINQF. COSQI stores in WK information needed for the fast Fourier computation of the discrete cosine and sine transforms and their inverses. A preliminary call must be made to COSQI before COSQB, COSQF, SINQB, and SINQF can be used. After this preliminary call, COSQI need only be recalled when  $n$  is modified.

**Programming.** COSQI employs the subroutines RFFTI and RFFTI1. These routines were written by Paul N. Swarztrauber (National Center for Atmospheric Research, Boulder, Colorado).

### CALL COSQB( $n$ , X, WK)

X is an array of dimension  $n$  or larger. On input it is assumed that X contains the data  $f(\theta_0), f(\theta_1), \dots, f(\theta_{n-1})$ . When COSQB is called,  $4na_j$  is computed and stored in  $X(j+1)$  for  $j = 0, 1, \dots, n-1$ .

WK is an array of dimension  $3n + 15$  or larger that is a work space for the routine. WK must be set up by the routine COSQI before COSQB can be used.

**Programming.** COSQB employs the subroutines COSB1, RFFTB, RFFTB1, RADB2, RADB3, RADB4, RADB5, and RADBG. These routines were written by Paul N. Swarztrauber (National Center for Atmospheric Research, Boulder, Colorado).

**CALL COSQF( $n, X, WK$ )**

$X$  is an array of dimension  $n$  or larger. On input it is assumed that  $X$  contains the data  $a_0, a_1, \dots, a_{n-1}$ . When COSQF is called,  $f(\theta_\nu)$  is computed and stored in  $X(\nu + 1)$  for  $\nu = 0, 1, \dots, n - 1$ .

$WK$  is an array of dimension  $3n + 15$  or larger that is a work space for the routine.  $WK$  must be set up by the routine COSQI before COSQF can be used.

**Example.** Assume that  $X$  contains the data  $f(\theta_0), \dots, f(\theta_{n-1})$ . When the statements

```
CALL COSQI(n,WK)
CALL COSQB(n,X,WK)
CALL COSQF(n,X,WK)
```

are called, COSQB stores  $4na_0, \dots, 4na_{n-1}$  in  $X$  and COSQF then sets  $X(\nu + 1) = 4nf(\theta_\nu)$  for  $\nu = 0, 1, \dots, n - 1$ . Thus, the terms of the original sequence  $X$  are multiplied by  $4n$ .

**Programming.** COSQF employs the subroutines COSQF1, RFFTF, RFFTF1, RADF2, RADF3, RADF4, RADF5, and RADFG. These routines were written by Paul N. Swarztrauber (National Center for Atmospheric Research, Boulder, Colorado).

**CALL SINQB( $n, X, WK$ )**

$X$  is an array of dimension  $n$  or larger. On input it is assumed that  $X$  contains the data  $f(\theta_0), \dots, f(\theta_{n-1})$ . When SINQB is called,  $4nb_j$  is computed and stored in  $X(j)$  for  $j = 1, \dots, n$ .

$WK$  is an array of dimension  $3n + 15$  or larger that is a work space for the routine.  $WK$  must be set up by the routine COSQI before SINQB can be used.

**Programming.** SINQB calls the subroutine COSQB. SINQB was written by Paul N. Swarztrauber (National Center for Atmospheric Research, Boulder, Colorado).

**CALL SINQF( $n, X, WK$ )**

$X$  is an array of dimension  $n$  or larger. On input it is assumed that  $X$  contains the data  $b_1, \dots, b_n$ . When SINQF is called,  $f(\theta_\nu)$  is computed and stored in  $X(\nu + 1)$  for  $\nu = 0, 1, \dots, n - 1$ .

$WK$  is an array of dimension  $3n + 15$  or larger that is a work space for the routine.  $WK$  must be set up by the routine COSQI before SINQF can be used.

**Example.** Assume that  $X$  contains the data  $b_1, \dots, b_n$ . When the statements

```
CALL COSQI(n,WK)
CALL SINQF(n,X,WK)
CALL SINQB(n,X,WK)
```



are called, SINQF stores  $f(\theta_0), \dots, f(\theta_{n-1})$  in  $X$  and SINQB then sets  $X(j) = 4nb_j$  for  $j = 1, \dots, n$ . Thus, the terms of the original sequence  $X$  are multiplied by  $4n$ .

**Programming.** SINQF calls the subroutine COSQF. The routine SINQF was written by Paul N. Swarztrauber (National Center for Atmospheric Research, Boulder, Colorado).

## RATIONAL MINIMAX APPROXIMATION OF FUNCTIONS

Let  $a < b$  and  $g(x)$  be a continuous nonvanishing function on the interval  $[a, b]$ . For any continuous function  $f(x)$ , let  $\|f\|$  denote the weighted norm  $\max\{|f(x)|/|g(x)| : a \leq x \leq b\}$ . Also let  $\phi(x)$  be a continuous strictly monotonic mapping on  $[a, b]$ . Then for any nonnegative integers  $\ell$  and  $m$ , the subroutine CHEBY is available for finding a rational function

$$R(x) = \frac{p_0 + p_1\phi(x) + \cdots + p_\ell\phi(x)^\ell}{q_0 + q_1\phi(x) + \cdots + q_m\phi(x)^m}$$

which minimizes  $\|R - f\|$ . The subroutine performs the calculations in double precision. It is assumed that the error curve  $\delta(x) = (R(x) - f(x))/g(x)$  satisfies  $|\delta(x_i)| = \|R - f\|$  at precisely  $\ell + m + 2$  critical points  $x_0 < x_1 < \cdots < x_n$  ( $n = \ell + m + 1$ ), and that  $\delta(x_{i+1}) = -\delta(x_i)$  for each  $i < n$ .

**CALL CHEBY**( $a, b, F, G, \text{PHI}, \epsilon, \text{ITER}, \text{MXITER}, \ell, m, P, Q,$   
 $\text{ERROR}, \text{IERR}, \text{WK}$ )

The arguments  $a$  and  $b$  are double precision real numbers.  $F, G$ , and  $\text{PHI}$  are functions whose arguments and values are double precision real numbers. The functions must be declared in the calling program to be of types DOUBLE PRECISION and EXTERNAL. The functions evaluate  $f(x)$ ,  $g(x)$ , and  $\phi(x)$  respectively.

The argument  $\epsilon$  is a double precision tolerance that is supplied by the user. If  $\lambda$  denotes the estimated value of  $\|R - f\|$ , then the routine converges when the error curve  $\delta(x)$  satisfies  $\lambda(1 - \epsilon) \leq |\delta(x_i)| \leq \lambda(1 + \epsilon)$  for each  $x_i$ . Thus  $\epsilon$  specifies the relative agreement that must be attained between  $\|f - R\|$  and the  $|\delta(x_i)|$ . Normally the setting  $\epsilon = 10^{-4}$  will give satisfactory results. It is required that  $0 < \epsilon < 10^{-2}$ .

The Remes-type algorithm designed by Cody, Fraser, and Hart is employed. This algorithm normally requires less than 20 iterations. The argument  $\text{MXITER}$  = the maximum number of iterations that may be performed. This argument is set by the user. The related argument  $\text{ITER}$  is an integer variable that is set by the routine. When CHEBY terminates,  $\text{ITER}$  will have for its value the number of iterations that were actually performed.

$P$  is a double precision array of dimension  $\ell + 1$ ,  $Q$  a double precision array of dimension  $m + 1$ ,  $\text{ERROR}$  a double precision variable, and  $\text{IERR}$  an integer variable. When CHEBY terminates, if the rational function approximation  $R(x)$  has been obtained then  $\text{IERR}$  is assigned the value 0 and  $\text{ERROR}$  is the estimated error  $\|R - f\|$ . The coefficient  $p_i$  of the numerator of  $R(x)$  is stored in  $P(i + 1)$  for  $i = 0, 1, \dots, \ell$ , and the coefficient  $q_j$  of the denominator is stored in  $Q(j + 1)$  for  $j = 0, 1, \dots, m$ . The coefficient  $q_0$  will always have the value 1.

Let  $k = \ell + m + 2$ . Then  $\text{WK}$  is a double precision array of dimension  $k(k + 5)$  or larger that is used for a work space.

**Error Return.** IERR is assigned one of the following values when the desired minimizing rational function  $R(x)$  is not obtained.

- IERR = 1 An input error was detected. Either  $\ell < 0$ ,  $m < 0$ ,  $\epsilon \leq 0$ ,  $\epsilon \geq 10^{-2}$ , or  $g(x) = 0$  for some point  $x$ .
- IERR = 2 MXITER iterations were performed. More iterations are needed to obtain  $R(x)$ .
- IERR = 3 The system of linear equations that define the coefficients  $p_i$  and  $q_j$  was found to be singular. This indicates that for the current values of  $\ell$  and  $m$ , the numerator and denominator of  $R(x)$  may have common factors.
- IERR = 4 A nonmonotonic sequence of critical points  $x_i$  was obtained. Modify  $\ell$  and/or  $m$ .
- IERR = 5 The value of the error curve  $\delta(x)$  at some critical point  $x_i$  appears to be too large. This indicates that  $R(x)$  may have poles, and that  $m$  (or possibly  $a$  or  $b$ ) may have to be modified.
- IERR = 6 CHEBY completely failed to find (or roughly approximate)  $R(x)$ . All information in  $P$ ,  $Q$ , and ERROR should be ignored.

If IERR = 2, 3, 4, or 5 then  $P$  and  $Q$  contain the coefficients of the most recent rational function approximation  $R(x)$  obtained, and ERROR is an estimate of the error  $\|R - f\|$  of the approximation.

**Remark.** The two most common weighting functions employed are  $g(x) = 1$  and  $g(x) = f(x)$ . If  $g(x) = 1$  then the absolute error is minimized in constructing  $R(x)$ . If  $g(x) = f(x)$  then the relative error is minimized.

**Programming.** CHEBY employs the subroutines CHEBY1, CERR, and DPSLV. These routines were written by A. H. Morris. CHEBY, CHEBY1, and CERR are slightly modified translations of the ALGOL 60 procedures Chebychev, lineq, del, and surmis given in the reference.

**Reference.** Cody, W. J., Fraser, W., and Hart, J. F., "Rational Chebychev Approximation using Linear Equations," *Numerische Mathematik* 12 (1968), pp. 242-251.

## L<sub>p</sub> APPROXIMATION OF FUNCTIONS

For any continuous real-valued function  $f(x)$  defined on the interval  $[a, b]$ , let  $\|f\|_p$  denote the  $L_p$  norm defined by

$$\begin{aligned}\|f\|_p &= \left( \int_a^b |f(x)|^p dx \right)^{1/p} & \text{if } 0 < p < \infty \\ \|f\|_p &= \max\{|f(x)| : a \leq x \leq b\} & \text{if } p = \infty.\end{aligned}$$

If  $p = \infty$  then the norm is also known as the *Chebyshev* norm. For any continuous function  $f$ ,  $0 < p \leq \infty$ , and  $\epsilon > 0$ , the subroutine ADAPT is available for finding a continuous piecewise polynomial function  $\phi$  that satisfies  $\|f - \phi\|_p \leq \epsilon$ .

**CALL ADAPT**( $F, a, b, \epsilon, k, \text{ERROR}, \text{XKNOTS}, C, \text{IND}, \mu, n, \ell, \text{ANORM},$   
 $\text{DX}, \text{MO}, m, \text{XBREAK}, \text{KDIFF}, \text{DLEFT}, \text{DRIGHT}$ )

It is assumed that the polynomials which form the approximation  $\phi$  are of degree  $\leq n$ . The argument  $n$  must satisfy  $1 \leq n \leq 19$ , and  $\text{IND}$  is a variable. When ADAPT is called, if there are no input errors and  $\phi$  is successfully constructed, then  $\text{IND}$  is set to 0, a sequence of points  $a = x_1 < \dots < x_{k-1} < x_k = b$  is selected, and  $\phi$  takes the form

$$\phi(x) = c_{i0} + c_{i1}(x - x_i) + \dots + c_{in}(x - x_i)^n \quad x_i \leq x \leq x_{i+1}$$

for  $i = 1, \dots, k-1$ . The points  $x_1, \dots, x_k$  are called the *knots* (or *nodes*) of  $\phi$ .

The argument  $\mu$  is the maximum number of polynomials that may be used in forming  $\phi$ .  $\text{ERROR}$  and  $k$  are variables,  $\text{XKNOTS}$  an array of dimension  $\mu + 1$  or larger, and  $C$  a 2-dimensional array of dimension  $\mu \times (n + 1)$ . ADAPT sets  $k$  to the number of knots that are generated. The knots  $x_1, \dots, x_k$  are stored in the  $\text{XKNOTS}$  array, and the coefficients  $c_{i0}, \dots, c_{in}$  are stored in  $C(i, 1), \dots, C(i, n + 1)$  for  $i = 1, \dots, k-1$ .  $\text{ERROR}$  is a rough estimate of the error  $\|f - \phi\|_p$ .

The argument  $\ell$  specifies the degree of smoothness that the approximation  $\phi$  must satisfy. It is assumed that  $0 \leq \ell < 10$  and  $n > 2\ell$ . If  $\ell = 0$  then it is only required that  $\phi$  be continuous on the interval  $[a, b]$ . Otherwise, if  $\ell \geq 1$  then it is assumed that  $f$  is of class  $C^\ell$  on  $[a, b]$  except at possibly a finite number of points (called *break points*), and it is required that  $\phi$  be of class  $C^\ell$  on  $[a, b]$  except possibly at the break points.

The argument  $m$  specifies the number of break points of  $f$ . It is assumed that  $m \leq 20$ . If  $m = 0$  then the arguments  $\text{XBREAK}, \text{KDIFF}, \text{DLEFT}$ , and  $\text{DRIGHT}$  can be ignored. Otherwise, if  $m \geq 1$  then it is assumed that  $\text{XBREAK}, \text{KDIFF}, \text{DLEFT}$ , and  $\text{DRIGHT}$  are arrays of dimension  $m$  or larger, and that

$\text{XBREAK}(i)$  = the  $i^{\text{th}}$  break point, call it  $u_i$ ,

$\text{KDIFF}(i)$  = the smallest integer  $\nu_i$  for which the  $\nu_i^{\text{th}}$  derivative of  $f$  does not exist or is not continuous at  $u_i$ ,

$\text{DLEFT}(i)$  = the value from the left of the  $\nu_i^{\text{th}}$  derivative at  $u_i$ , and

$\text{DRIGHT}(i)$  = the value from the right of the  $\nu_i^{\text{th}}$  derivative at  $u_i$

for  $i = 1, \dots, m$ . It is also assumed that  $a \leq u_1 < \dots < u_m \leq b$  and  $n > 2\nu_i$  for each  $\nu_i$ .

$F$  is the name of a user defined function that has the value  $F(x, D) = f(x)$  for  $a \leq x \leq b$ . If  $\ell = 0$  then  $D$  can be ignored. (However,  $D$  must still be given as an argument of  $F$ .) Otherwise, if  $\ell \geq 1$  then  $D$  is an array of dimension greater than or equal to  $\ell$ . For any  $x$  not a break point in XBREAK, the user must set  $D(j) =$  the  $j^{\text{th}}$  derivative of  $f$  at  $x$  for  $j \leq \ell$ . However, if  $x = \text{XBREAK}(i)$  then the user need only set  $D(j) =$  the  $j^{\text{th}}$  derivative of  $f$  at  $x$  for  $j < \text{KDIFF}(i)$ . The function  $F$  must be declared in the calling program to be of type EXTERNAL.

The argument DX specifies the maximum distance to be permitted between the knots  $x_i$ , and the argument ANORM specifies the norm to be used. Set

ANORM =  $\pm 1.0$  for  $L_1$  approximation.

ANORM =  $\pm 2.0$  for  $L_2$  (least squares) approximation.

ANORM =  $3.0$  for  $L_\infty$  (minimax) approximation.

ANORM =  $-p$  for  $L_p$  ( $0 < p < \infty$ ) approximation.

Before considering the argument MO, one should be briefly acquainted with how ADAPT operates. ADAPT employs the following procedure to construct  $\phi$ .

- (1) Set  $I = [a, b]$  and  $k = 1$ . Let  $a$  be the first knot of  $\phi$ .
- (2) If the interior of  $I$  contains no break points then go to (3). Otherwise, if  $I = [c, d]$  then partition  $I$  into the subintervals  $[c, u]$  and  $[u, d]$  where  $u$  is the smallest break point greater than  $c$ . Stack the right subinterval  $[u, d]$  and reset  $I$  to  $[c, u]$ .
- (3) Construct a polynomial  $\phi_I$  on  $I$  using Hermite interpolation. If the length of the interval  $I$  is  $\leq \text{DX}$  and  $\phi_I$  satisfactorily approximates  $f$  on  $I$ , then go to (4). Otherwise go to (5).
- (4) Set  $k$  to be  $k + 1$ . Let  $\phi_I$  be the  $(k - 1)^{\text{st}}$  polynomial forming  $\phi$  and let the right end point of  $I$  be the  $k^{\text{th}}$  knot of  $\phi$ . If the interval stack is empty then the procedure is finished. Otherwise, obtain from the stack the next interval  $I$  to be considered and return to (2).
- (5) The polynomial  $\phi_I$  cannot be used. Partition  $I$  into halves, stack the right subinterval and reset  $I$  to be the left subinterval. Then go to (3).

The argument MO specifies the accuracy criterion that the approximation  $\phi$  is to satisfy on a subinterval  $I = [c, d]$  of  $[a, b]$ . It is assumed that  $\text{MO} = 0, 1, 2$ . If the  $L_\infty$  norm is used then MO is ignored<sup>1</sup> and  $\phi$  is required to satisfy  $|f(x) - \phi(x)| \leq \epsilon$  for  $c \leq x \leq d$ . Otherwise, if the  $L_p$  ( $0 < p < \infty$ ) norm is used then  $\phi$  is required to satisfy:

$$\int_c^d |f(x) - \phi(x)|^p dx \leq \frac{d-c}{b-a} \epsilon^p \quad \text{for MO} = 0$$

$$\int_c^d |f(x) - \phi(x)|^p dx \leq \epsilon^p \quad \text{for MO} = 2$$

The setting  $\text{MO} = 0$ , which is the most commonly used setting, requires the total error  $\|f - \phi\|_p \leq \epsilon$ . The alternate setting  $\text{MO} = 2$  employs  $\epsilon$  to control local accuracy. If  $\phi$  consists of  $k - 1$  polynomials then the total error  $\|f - \phi\|_p \leq (k - 1)^{1/p} \epsilon$ . This setting can

<sup>1</sup> However, it is still required that  $\text{MO} = 0, 1, 2$ .

be useful when  $f$  is rough. A (heuristic) compromise strategy is provided when  $MO = 1$ . At each step in the formation of  $\phi$ , the  $MO = 1$  strategy estimates the total number of subintervals that will finally be needed and adjusts the error requirement for the subinterval  $I$  accordingly. This strategy attempts to keep the total error to a minimum while relaxing the local accuracy criterion demanded by the  $MO = 0$  setting.

**Remarks.**  $IND$ ,  $k$ ,  $\mu$ ,  $n$ ,  $\ell$ ,  $MO$ ,  $m$ , and  $KDIFF$  are integer arguments. All other arguments (including  $F$ ) are double precision arguments.

**Error Return.** ADAPT assigns  $IND$  one of the following values:

- $IND = 0$  The approximation was successfully constructed.
- $IND = -1$  Either  $a \geq b$  or one of the arguments  $\epsilon$ ,  $n$ ,  $\ell$ ,  $ANORM$ ,  $MO$ ,  $m$  is assigned an incorrect value.
- $IND = -2$   $[a, b]$  is too small an interval.
- $IND = -3$   $DX$  is less than  $(b - a)/\mu$ . Since only  $\mu$  subintervals can be used and each subinterval must be of length  $\leq DX$ , the interval  $[a, b]$  cannot be covered. Make  $DX$  or  $\mu$  larger.
- $IND = -4$  The restriction  $a \leq u_1 < \dots < u_m \leq b$  on the break points is violated.
- $IND = -5$  Either  $KDIFF(i) < 0$  or  $KDIFF(i) > (n - 1)/2$  for some  $i$ .
- $IND = 1$  ADAPT selected  $\mu + 1$  knots. More knots are needed to complete the problem.
- $IND = 2$  A subinterval  $I = [c, d]$  must be partitioned into subintervals  $[c, u]$  and  $[u, d]$  where  $u$  is a break point. However, this cannot be done either because the interval stack is full, or partitioning will produce too small an interval. (The stack can hold only 50 subintervals).
- $IND = 3$  A subinterval must be partitioned because its length is greater than  $DX$ . However, this cannot be done since the interval stack is full.
- $IND = 4$  A subinterval must be partitioned so that the accuracy criterion can be satisfied. However, this cannot be done either because the stack is full, or partitioning will produce too small an interval.

If an input error is detected (i.e., if  $IND < 0$ ) then no computation is performed. Otherwise, if  $IND \geq 0$  then when ADAPT terminates  $k$  = the number of knots generated,  $XKNOTS$  contains the knots,  $C$  contains the coefficients of the polynomials generated, and  $ERROR$  contains the error estimate for  $f - \phi$  over the interval covered.

**Remarks.**

- (1) If the  $L_\infty$  norm is used then  $\epsilon$  controls absolute accuracy, not relative accuracy. This should be kept in mind when  $\epsilon$  is to be set for any  $L_p$  norm.
- (2) ADAPT requires more time when  $\ell \geq 2$  than when  $\ell = 0$  or 1. However, the choice of the norm normally has little effect on the efficiency of the routine.
- (3) ADAPT can yield excellent results even when the derivatives of  $f$  have singularities. The one major exception is when the first derivative of  $f$  is not bounded. Then the routine can be expected to fail.

**Example.** The following code can be used for approximating  $f(x) = e^x$  on the interval  $[0, 1]$ .

```
DOUBLE PRECISION F, A, B, EPS, ERROR, ANORM, DX
DOUBLE PRECISION XKNOTS (11), C(10,20)
INTEGER KDIFF(1)
DOUBLE PRECISION XBREAK(1),DLEFT(1),DRIGHT(1)
EXTERNAL F
DATA MAX, A, B, DX/10, 0.D0, 1.D0, 1.D0/
N = 8
L = 1
EPS = 1.D-12
ANORM = 3.D0
CALL ADAPT(F,A,B,EPS,K,ERROR,XKNOTS,C,IND,
* MAX,N,L,ANORM,DX,0,0,XBREAK,KDIFF,DLEFT,DRIGHT)
```

Here  $F$  may be defined by:

```
DOUBLE PRECISION FUNCTION F(X,D)
DOUBLE PRECISION X,D(1)
F = DEXP(X)
D(1) = F
RETURN
END
```

In the ADAPT statement XBREAK, KDIFF, DLEFT, and DRIGHT are ignored since  $m = 0$ .

**Programming.** ADAPT employs the subroutines ADAPT1, ADSET, ADTAKE, ADCOMP, NEWTON, ADCHK, ADPUT, ADTRAN and functions ERRINT, POLYDD. These routines exchange information in labeled common blocks. The block names are INPUTZ, RESULTZ, KONTRL, and COMDIF. The routines were written by John R. Rice (Purdue University) and modified by A. H. Morris. The function DPMPAR is also used.

#### References.

- (1) Rice, J. R., "Algorithm 525. ADAPT, Adaptive Smooth Curve Fitting," *ACM Trans. Math Software* 4 (1978), pp. 82-94.
- (2) \_\_\_\_\_, "Adaptive Approximation," *J. Approx. Theory* 16 (1976), pp. 329-337.

## CALCULATION OF THE TAYLOR SERIES OF A COMPLEX ANALYTIC FUNCTION

Let  $f(z) = \sum_{n \geq 0} a_n(z - z_0)^n$  denote the Taylor series of an analytic function  $f$  around a point  $z_0$ . Then the subroutines CPSC and DCPSC are available for obtaining the coefficients  $a_n$  of the series. CPSC obtains single precision results and DCPSC obtains double precision results.

### CALL CPSC( $f, z_0, n, \text{IND}, \epsilon, R, A, \text{ERR}$ )

It is assumed that  $f(z)$  is a user defined function whose arguments and values are complex numbers. The argument  $f$  must be declared in the calling program to be of types COMPLEX and EXTERNAL.

The argument  $z_0$  is complex,  $n$  is an integer where  $1 \leq n \leq 51$ , and  $A$  is a complex array of dimension  $n$  or larger. IND may be any integer. If  $\text{IND} = 0$  then  $a_j$  is computed and stored in  $A(j+1)$  for  $j = 0, 1, \dots, n-1$ . Otherwise, if  $\text{IND} \neq 0$  then  $f(z_0)$  and the derivatives  $f'(z_0), \dots, f^{(n-1)}(z_0)$  are computed and stored in  $A$ .

The argument  $\epsilon$  specifies the relative accuracy of  $f$ . If it is estimated that  $f$  produces results accurate to  $k$  significant decimal digits then one may set  $\epsilon = 10^{-k}$ . It is assumed that  $\epsilon \geq 0$ . If  $\epsilon = 0$  then the results of  $f$  are assumed to be correct to machine precision.

When CPSC is called,  $f(z)$  is evaluated on circles of various radii around the point  $z_0$ .  $R$  is a real variable. On input,  $R$  is the radius of the first circle on which  $f(z)$  is to be evaluated. After using this radius, the radius is repeatedly modified (first by factors of 2 or 1/2) until a suitable final radius  $r_c$  is obtained for deriving the values of the coefficients  $a_j$ . This radius, whose value depends on  $\epsilon$ , is called the computational radius of the series  $\sum a_j(z - z_0)^j$ . When the routine terminates,  $R$  is assigned the value  $r_c$ .

ERR is a real array of dimension  $n$  or larger. On output,  $\text{ERR}(j)$  is the estimated absolute error of  $A(j)$  for  $j = 1, \dots, n$ .

**Usage.** Given a radius  $R$ ,  $f(z)$  is evaluated on  $k$  equidistant points on the circle of radius  $R$  around  $z_0$  where

$$\begin{aligned} k &= 8 \text{ when } 1 \leq n \leq 6, \\ k &= 16 \text{ when } 7 \leq n \leq 12, \\ k &= 32 \text{ when } 13 \leq n \leq 25, \\ k &= 64 \text{ when } 26 \leq n \leq 51. \end{aligned}$$

It is assumed that  $f(z)$  has at least one nonzero coefficient  $a_j$  among the first  $k/2$  coefficients, and at least one nonzero coefficient among the next  $k/2$  coefficients. Thus, the routine should not be used to obtain coefficients of a low degree polynomial such as  $f(z) = 1 - z^2$ . In such cases, the results will normally be incorrect.

In general, the selection of the radius  $R$  of the first circle on which  $f(z)$  is evaluated is not bothersome. A randomly selected value of moderate size, such as  $R = 6.2738$ , almost



always suffices. No difficulties normally arise when the initial radius  $R$  is greater than the radius of convergence of the series  $\sum a_j(z - z_0)^j$ . However, difficulties do arise when

- (1) the routine attempts to evaluate  $f$  too close to (or exactly at) a singularity,
- (2)  $f$  has a Taylor series expansion which contains one or more extremely large isolated terms (e.g.,  $f(z) = 10^8 + \sin z$ ),
- (3)  $f$  has a branch point near  $z_0$ , or
- (4) The initial radius  $R$  is too far from the computational radius  $r_c$  (see the error return section below).

The risk of (1) occurring is minimized by the random selection of an initial radius  $R$ . For (2) and (3), a severe loss of accuracy can occur when a large number of coefficients are requested. In these cases, any loss of accuracy is reported by the ERR array.

**Error return.**  $A(j)$  is assigned the value 0 and  $ERR(j)$  is set to  $10^{10}$  for  $j = 1, \dots, n$  when the initial radius  $R$  differs from  $r_c$  by a factor of 30000 or greater.

**Programming.** CPSC was written by Bengt Fornberg (California Institute of Technology) and modified by A.H. Morris. CPSC employs the function SPMPAR.

#### References.

- (1) Fornberg, B., "Numerical Differentiation of Analytic Functions," *ACM Trans. Math Software* 7, 1981, pp. 512-526.
- (2) ———, "Algorithm 579. CPSC: Complex Power Series Coefficients," *ACM Trans. Math Software* 7, 1981, pp. 542-547.

#### CALL DCPSC( $F, x_0, y_0, n, IND, \epsilon, R, AR, AI, ERR$ )

$F$  is the name of a user defined subroutine that has the format:

CALL  $F(x, y, u, v)$

This subroutine is used for evaluating  $f(z)$  at point  $z$ . The arguments  $x$  and  $y$  are the real and imaginary parts of  $z$ , and  $u$  and  $v$  are the real and imaginary parts of  $f(z)$ . The arguments  $x$  and  $y$  have double precision values, and  $u$  and  $v$  are double precision variables.  $F$  must be declared in the calling program to be of type EXTERNAL.

The arguments  $x_0$  and  $y_0$ , which have double precision values, are the real and imaginary parts of  $z_0$ . The argument  $n$  is an integer where  $1 \leq n \leq 51$ , and  $AR$  and  $AI$  are double precision arrays of dimension  $n$  or larger.  $IND$  may be any integer. If  $IND = 0$  then the real and imaginary parts of  $a_j$  are stored in  $AR(j+1)$  and  $AI(j+1)$  for  $j = 0, 1, \dots, n-1$ . Otherwise, if  $IND \neq 0$  then the real and imaginary parts of  $f(z_0)$  and the derivatives  $f'(z_0), \dots, f^{(n-1)}(z_0)$  are stored in  $AR$  and  $AI$  respectively.

The argument  $\epsilon$ , which has double precision values, specifies the relative accuracy of the subroutine  $F$ . If it is estimated that  $F$  produces results accurate to  $k$  decimal digits then one may set  $\epsilon = 10^{-k}$ . It is assumed that  $\epsilon \geq 0$ . If  $\epsilon = 0$  then the results of  $F$  are assumed to be correct to machine precision.

When DCPSC is called,  $f(z)$  is evaluated on circles of various radii around the point  $z_0$ .  $R$  is a double precision variable. On input,  $R$  is the radius of the first circle on which

$f(z)$  is to be evaluated. After using this radius, the radius is repeatedly modified (first by factors of 2 or 1/2) until a suitable final radius  $r_c$  is obtained for deriving the values of the coefficients  $a_j$ . This radius, whose value depends on  $\epsilon$ , is called the computational radius of the series  $\sum a_j(z - z_0)^j$ . When the routine terminates  $R$  is assigned the value  $r_c$ .

ERR is a double precision array of dimension  $n$  or larger. On output,  $ERR(j)$  is the estimated absolute error of the complex value stored in  $AR(j)$  and  $AI(j)$  for  $j = 1, \dots, n$ .

**Usage.** DCPSC is used in the same manner as CPSC. See the usage section for CPSC.

**Error return.**  $AR(j)$  and  $AI(j)$  are assigned the value 0 and  $ERR(j)$  is set to  $10^{10}$  for  $j = 1, \dots, n$  when the initial radius  $R$  differs from  $r_c$  by a factor of 30000 or greater.

**Programming.** DCPSC is an adaptation by A.H. Morris of the subroutine CPSC, written by Bengt Fornberg (California Institute of Technology). DCPSC employs the function DPMPAR.

#### References.

- (1) Fornberg, B., "Numerical Differentiation of Analytic Functions," *ACM Trans. Math Software* 7, 1981, pp. 512-526.
- (2) \_\_\_\_\_, "Algorithm 579. CPSC: Complex Power Series Coefficients," *ACM Trans. Math Software* 7, 1981, pp. 542-547.

## LINEAR INTERPOLATION

Let  $a$  be a real number and  $(x_1, y_1), \dots, (x_n, y_n)$  a sequence of points. The following function performs a linear interpolation at point  $a$ .

**TRP**( $a, n, X, Y$ )

It is assumed that  $n \geq 2$  and  $x_1 < \dots < x_n$ .  $X$  and  $Y$  are arrays containing the abscissas  $x_1, \dots, x_n$  and ordinates  $y_1, \dots, y_n$  respectively.  $\text{TRP}(a, n, X, Y) = b$  where  $b$  is the value of the interpolation at  $a$ .

**Programmer.** A. H. Morris

## LAGRANGE INTERPOLATION

Let  $\{(x_i, y_i) : i = 1, \dots, n\}$  be a set of  $n \geq 2$  points where  $x_1 < \dots < x_n$ ,  $m$  be an integer where  $2 \leq m \leq n$ , and  $\bar{x}_1, \dots, \bar{x}_k$  be  $k \geq 1$  points at which  $m$  point Lagrange interpolation is to be performed. The subroutine LTRP is available for performing this interpolation.

**CALL LTRP( $m, X, Y, n, XI, YI, k, T, IERR$ )**

$X$  is an array containing  $x_1, \dots, x_n$ ,  $Y$  an array containing  $y_1, \dots, y_n$ ,  $XI$  an array containing  $\bar{x}_1, \dots, \bar{x}_k$ , and  $YI$  an array of dimension  $k$  or larger. When LTRP is called, if no input errors are detected then interpolation is performed at each  $\bar{x}_j$  and the result stored in  $YI(j)$  for  $j = 1, \dots, k$ .

$T$  is an array of dimension  $m$  or larger. The array is used as a temporary storage area by the routine.

**Error Return.** IERR is a variable that is set by the routine. If no input errors are detected then IERR is assigned the value 0. Otherwise, IERR is assigned one of the following values:

IERR = 1 if  $m < 2$ .

IERR = 2 if  $m > n$ .

IERR = 3 if  $k < 1$ .

When an error is detected LTRP immediately terminates.

**Algorithm.** If  $\bar{x}_j = (x_i + x_{i+m})/2$  for some  $i$ , then  $(x_i, y_i), \dots, (x_{i+m-1}, y_{i+m-1})$  are the  $m$  data points used in the Lagrange interpolation at  $\bar{x}_j$ . Otherwise, the data points selected for the interpolation are those  $m$  points  $(x_i, y_i)$  whose abscissas are closest to  $\bar{x}_j$ .

**Linear Interpolation.** For  $m = 2$ , if the abscissae  $x_i$  are not equally spaced then LTRP can produce different results than the linear interpolating function TRP. If  $\bar{x}_j$  lies in the interval  $[x_i, x_{i+1})$  then TRP always uses the data points  $(x_i, y_i)$  and  $(x_{i+1}, y_{i+1})$  to find the interpolated value at  $\bar{x}_j$ . However, the operation of LTRP is somewhat different. For example, if the point  $\bar{x}_j$  in  $[x_i, x_{i+1})$  is closer to  $x_{i-1}$  than to  $x_{i+1}$ , then  $(x_{i-1}, y_{i-1})$  and  $(x_i, y_i)$  will be the data points employed in the interpolation. Thus, TRP will normally be the procedure that will be desired for linear interpolation.

**Programming.** Developed by A. H. Morris. The portion of the code for finding the subinterval containing  $\bar{x}_j$  was written by Rondall E. Jones (Sandia Laboratories).

## HERMITE INTERPOLATION

Let  $x_1, \dots, x_k$  be  $k \geq 1$  distinct points. For each  $x_i$  assume that we are given  $n_i \geq 1$  values  $y_i, y'_i, \dots, y_i^{(n_i-1)}$ . If  $n = n_1 + \dots + n_k$ , then there exists a unique polynomial of degree  $n - 1$  which satisfies

$$\begin{aligned} p(x_i) &= y_i \\ p'(x_i) &= y'_i \\ &\vdots \\ p^{n_i-1}(x_i) &= y_i^{(n_i-1)} \end{aligned}$$

for each  $i = 1, \dots, k$ . The subroutine HTRP is available for obtaining this polynomial.

**CALL HTRP( $n, X, Y, A, WK, IERR$ )**

$X$  and  $Y$  are arrays of dimension  $n$  containing the following information:  $X(j) = x_1$  for  $j = 1, \dots, n_1$  and  $Y(1), \dots, Y(n_1)$  contain the values  $y_1, y'_1, \dots, y_1^{(n_1-1)}$ . For  $i = 2, \dots, k$  let  $m_i = n_1 + \dots + n_{i-1}$ . Then  $X(m_i + j) = x_i$  for  $j = 1, \dots, n_i$  and  $Y(m_i + 1), \dots, Y(m_i + n_i)$  contain the values  $y_i, y'_i, \dots, y_i^{(n_i-1)}$ .

$A$  is an array of dimension  $n$  and  $IERR$  an integer variable. When HTRP is called, if no errors are detected then  $IERR$  is assigned the value 0 and the coefficients  $a_j$  of the polynomial  $p(x) = a_0 + \sum_{j=1}^{n-1} a_j(x - X(1)) \cdots (x - X(j))$  are computed and stored in  $A(j+1)$  for  $j = 0, 1, \dots, n-1$ .

$WK$  is an array of dimension  $n$  or larger that is a work space for the routine.

**Error Return.** If an error is detected then  $IERR$  is assigned one of the following values:

$IERR = 1$  The argument  $n$  is not positive.

$IERR = 2$  There exists integers  $i$  and  $\ell$  for which  $X(i) = X(\ell)$  but  $X(i) \neq X(j)$  for some  $j$  where  $i < j < \ell$ . In this case, the values  $i$  and  $\ell$  are stored (in floating point form) in  $WK(1)$  and  $WK(2)$ .

When an error is detected, the routine immediately terminates.

**Example.** If  $p(0) = 2$ ,  $p(-1) = 1$ , and  $p'(-1) = 2$  where  $x_1 = 0$  and  $x_2 = -1$ , then HTRP stores 2, 1, -1 in  $A$ . Hence,  $p(x) = 2 + x - x(x+1)$  is the desired polynomial.

**Remark.** The Newton representation  $a_0 + \sum_{j=1}^{n-1} a_j(x - X(1)) \cdots (x - X(j))$  of the polynomial  $p(x)$  can be converted to the Taylor series representation  $\sum_{j=0}^{n-1} c_j(x - \alpha)^j$  by the subroutine PCOEFF.

**Programmer.** A. H. Morris

## CONVERSION OF REAL POLYNOMIALS FROM NEWTON TO TAYLOR SERIES FORM

For  $n \geq 1$  let  $p(x) = a_0 + \sum_{j=1}^{n-1} a_j(x - x_1) \cdots (x - x_j)$ . Then for any real number  $\alpha$ , the subroutine PCOEFF is available for converting the polynomial  $p(x)$  to the Taylor series form  $\sum_{j=0}^{n-1} c_j(x - \alpha)^j$ .

**CALL PCOEFF**( $\alpha, n, X, A, C, T$ )

$X$  is a single precision real array containing  $x_1, \dots, x_{n-1}$ ,  $A$  a single precision real array containing  $a_0, a_1, \dots, a_{n-1}$  where  $a_j$  is stored in  $A(j+1)$  for  $j = 0, 1, \dots, n-1$ , and  $C$  a single precision real array of dimension  $n$  or larger. When PCOEFF is called then the coefficients  $c_j$  of the Taylor series representation are computed and stored in  $C(j+1)$  for  $j = 0, 1, \dots, n-1$ .

$T$  is a *double precision* array of dimension  $n$  or larger. The array is a work space for the routine. (The conversion of the coefficients is done in double precision.)

**Note.**  $A$  and  $C$  may reference the same storage area, in which case the results  $c_j$  will overwrite the input data  $a_j$ .

**Programmer.** A. H. Morris

## LEAST SQUARES POLYNOMIAL FIT

Let  $\{(x_i, y_i) : i = 1, \dots, m\}$  be a set of  $m \geq 2$  points where  $x_i \neq x_j$  for  $i \neq j$ . Then for any positive integer  $n$  where  $n < m$ , the subroutine PFIT is available for obtaining the (unique)  $n^{\text{th}}$  degree polynomial  $p(x) = \sum_{j=0}^n a_j x^j$  which minimizes  $\sum_{i=1}^m (p(x_i) - y_i)^2$ .

**CALL PFIT**( $n, m, X, Y, A, \text{RNORM}, \text{PHI}, \text{WK}, \text{IERR}$ )

$X$  is an array containing  $x_1, \dots, x_m$ ,  $Y$  an array containing  $y_1, \dots, y_m$ , and  $A$  an array of dimension  $n+1$  or larger.  $\text{RNORM}$  and  $\text{IERR}$  are variables. When PFIT is called, if no input errors are detected then  $\text{IERR}$  is set to 0, the coefficients  $a_j$  of  $p(x)$  are stored in  $A(j+1)$  for  $j = 0, 1, \dots, n$ , and  $\text{RNORM}$  is assigned the value  $\sqrt{\sum_{i=1}^m (p(x_i) - y_i)^2}$ .

$\text{PHI}$  is an array of dimension  $2(n+1)$  or larger, and  $\text{WK}$  is an array of dimension  $4m$  or larger.  $\text{PHI}$  and  $\text{WK}$  are work spaces for the routine.

**Error Return.**  $\text{IERR} = 1$  if  $n < 1$  or  $n \geq m$ .

**Algorithm.** The abscissas  $x_i$  are first mapped into values in the interval  $[-1, 1]$ . Then the Forsythe procedure is used.

**Programmer.** A. H. Morris

## WEIGHTED LEAST SQUARES POLYNOMIAL FIT

Let  $\{(x_i, y_i) : i = 1, \dots, m\}$  be a set of  $m \geq 2$  points where  $x_i \neq x_j$  for  $i \neq j$ , and let  $w_i \geq 0$  ( $i = 1, \dots, m$ ) be weights. It is assumed that  $m_w \geq 2$  where  $m_w$  is the number of nonzero weights. For any positive integer  $n$  where  $n < m_w$ , the subroutine WPFIT is available for finding the (unique)  $n^{\text{th}}$  degree polynomial  $p(x) = \sum_{j=0}^n a_j x^j$  which minimizes  $\sum_{i=1}^m w_i (p(x_i) - y_i)^2$ .

**CALL WPFIT**( $n, m, X, Y, W, A, \text{RNORM}, \text{PHI}, \text{WK}, \text{IERR}$ )

$X$  is an array containing  $x_1, \dots, x_m$ ,  $Y$  an array containing  $y_1, \dots, y_m$ ,  $W$  an array containing  $w_1, \dots, w_m$ , and  $A$  an array of dimension  $n + 1$  or larger.  $\text{RNORM}$  and  $\text{IERR}$  are variables. When WPFIT is called, if no input errors are detected then  $\text{IERR}$  is set to 0, the coefficients  $a_j$  of  $p(x)$  are computed and stored in  $A(j + 1)$  for  $j = 0, 1, \dots, n$ , and  $\text{RNORM}$  is assigned the value  $\sqrt{\sum_i w_i (p(x_i) - y_i)^2}$ .

$\text{PHI}$  is an array of dimension  $2(n + 1)$  or larger, and  $\text{WK}$  is an array of dimension  $4m$  or larger.  $\text{PHI}$  and  $\text{WK}$  are work spaces for the routine.

**Error Return.**  $\text{IERR} = 1$  if  $n < 1$  or  $n \geq m_w$ , and  $\text{IERR} = 3$  if some  $w_i$  is negative.

**Algorithm.** The abscissas  $x_i$  corresponding to the positive weights are first mapped into values in the interval  $[-1, 1]$ . Then the Forsythe procedure is used.

**Programmer.** A. H. Morris



## CUBIC SPLINE INTERPOLATION

Given  $x_1 < \dots < x_n$ . A function  $f(x)$  is a *cubic spline having the nodes(knots)*  $x_1, \dots, x_n$  if  $f$  is a polynomial of degree  $\leq 3$  on the interval  $[x_i, x_{i+1}]$  for  $i = 1, \dots, n-1$ , and the first and second derivatives  $f'(x)$  and  $f''(x)$  exist and are continuous for all  $x$ . If  $f_i$  denotes the polynomial for the interval  $[x_i, x_{i+1}]$  then  $f_i$  has the form:

$$f_i(x) = y_i + a_i(x - x_i) + b_i(x - x_i)^2 + c_i(x - x_i)^3.$$

Consequently, the function  $f$  is obtained by fitting the polynomials  $f_1, \dots, f_{n-1}$  together at the points  $x_2, \dots, x_{n-1}$ . For  $x < x_1$   $f(x) = f_1(x)$ , and for  $x > x_n$   $f(x) = f_{n-1}(x)$ . Also  $f(x_i) = y_i$  for  $i = 1, \dots, n-1$ . Hence, if  $f(x_n) = y_n$  then  $f$  interpolates the points  $(x_i, y_i)$  for  $i = 1, \dots, n$ .

Assume now that the ordinates  $y_1, \dots, y_n$  are given. Then there exist an infinitude of splines with nodes  $x_1, \dots, x_n$  that interpolate the points  $(x_i, y_i)$ . In general, two independent conditions must be imposed to uniquely specify the interpolating spline  $f$  which is of interest. Frequently,  $f$  is restricted on the first interval  $[x_1, x_2]$  by requiring that  $f'(x_1)$  or  $f''(x_1)$  has a given value, or that  $f$  is continuous at  $x_2$ . Also,  $f$  is restricted on the last interval  $[x_{n-1}, x_n]$  by requiring that  $f'(x_n)$  or  $f''(x_n)$  has a given value, or that  $f$  is continuous at  $x_{n-1}$ . The subroutine CBSPL is available for obtaining the spline when these restrictions are employed. Alternatively,  $f$  may be required to satisfy the conditions

$$\begin{aligned} f''(x_1) &= \alpha f''(x_2) + \beta & |\alpha| < 1 \\ f''(x_n) &= \bar{\alpha} f''(x_{n-1}) + \bar{\beta} & |\bar{\alpha}| < 1 \end{aligned}$$

when  $n \geq 4$ . Then the subroutine SPLIFT is available for obtaining the spline.

**CALL CBSPL(X,Y,A,B,C,n,i<sub>1</sub>,i<sub>2</sub>,w<sub>1</sub>,w<sub>2</sub>,IERR)**

$X$  and  $Y$  are arrays containing the abscissas  $x_1, \dots, x_n$  and ordinates  $y_1, \dots, y_n$ . It is assumed that  $x_1 < \dots < x_n$  and  $n \geq 3$ .  $A, B$  and  $C$  are arrays of dimension  $n$  or larger, and IERR an integer variable. When CBSPL is called, if no input errors are detected then IERR is set to 0. Also, the coefficients  $a_i, b_i, c_i$  ( $i = 1, \dots, n-1$ ) of the interpolating spline  $f(x)$  are stored in  $A, B, C$ , and  $A(n)$  is set to  $f'(x_n)$ .

The arguments  $i_1, i_2, w_1, w_2$  specify the conditions that the spline  $f(x)$  must satisfy. It is assumed that  $i_1$  and  $i_2$  have the values 0,1,2 where:

$i_1 = 0$ $f$ is continuous at $x_2$ . $i_1 = 1$ $f'(x_1)$ has the value $w_1$ . $i_1 = 2$ $f''(x_1)$ has the value $w_1$ .	$i_2 = 0$ $f$ is continuous at $x_{n-1}$ . $i_2 = 1$ $f'(x_n)$ has the value $w_n$ . $i_2 = 2$ $f''(x_n)$ has the value $w_n$ .
-----------------------------------------------------------------------------------------------------------------------------------	---------------------------------------------------------------------------------------------------------------------------------------

If  $i_1 = 0$  then the argument  $w_1$  is ignored, and if  $i_2 = 0$  then  $w_2$  is ignored.

**Error Return.** IERR = 1 if  $n < 3$  and IERR = 2 if  $x_i \geq x_{i+1}$  for some  $i$ .

**Remarks.**

- (1)  $B(n)$  and  $C(n)$  are used for temporary storage.
- (2) If  $i_1 = i_2 = 0$  and  $n = 3$ , then it is also assumed that  $f'(x_2) + f'(x_3) = 2\left(\frac{y_3 - y_2}{x_3 - x_2}\right)$ .
- (3) After  $A, B$  and  $C$  have been obtained, then SCOMP or SCOMP1 may be used to evaluate the spline at any point  $x$ . SEVAL or SEVAL1 may be used if derivatives are also desired.

**Programming.** CBSPL is an adaptation by A. H. Morris of the subroutine CUBSPL, written by Carl de Boor (University of Wisconsin).

**Reference.** de Boor, Carl, *A Practical Guide to Splines*, Springer-Verlag, 1978.

**CALL SPLIFT**( $X, Y, DY, DDY, n, W, IERR, MO, \alpha, \beta, \bar{\alpha}, \bar{\beta}$ )

$X$  and  $Y$  are arrays containing the abscissas  $x_1, \dots, x_n$  and ordinates  $y_1, \dots, y_n$ . It is assumed that  $x_1 < \dots < x_n$  and  $n \geq 4$ .  $DY$  and  $DDY$  are arrays of dimension  $n$  or larger, and  $IERR$  an integer variable. When SPLIFT is called, if there are no input errors then  $IERR$  is assigned the value 0, the first derivatives  $f'(x_1), \dots, f'(x_n)$  are computed and stored in  $DY(1), \dots, DY(n)$ , and the second derivatives  $f''(x_1), \dots, f''(x_n)$  are computed and stored in  $DDY(1), \dots, DDY(n)$ .

$W$  is an array of dimension  $3n$  or larger that is used for a storage area. On the first call to SPLIFT the argument  $MO$  must be set to 0. When SPLIFT is initially called, certain calculations which depend only on the values of  $\alpha, \bar{\alpha}$ , and  $x_1, \dots, x_n$  are performed and the results stored in  $W$ . On subsequent calls to SPLIFT, if only values of  $\beta, \bar{\beta}$  or  $y_1, \dots, y_n$  are modified, then the information in  $W$  need not be recomputed. Set  $MO = 1$  and the information in  $W$  will be reused.

**Error Return** If there is an input error then  $IERR$  is set as follows:

$IERR = 1$  if  $|\alpha| \geq 1$  or  $|\bar{\alpha}| \geq 1$ .

$IERR = 2$  if  $n < 4$ .

$IERR = 3$  if the restriction  $x_1 < \dots < x_n$  is not satisfied.

**Remarks.**

- (1) After  $DY$  and  $DDY$  have been obtained, then SCOMP1 or SCOMP2 may be used to evaluate the spline at any point  $x$ . SEVAL1 or SEVAL2 may be used if derivatives are also desired.
- (2) Given the values  $y'_1$  and  $y'_n$ . Then there exists a unique interpolating cubic spline  $f$  that satisfies  $f'(x_1) = y'_1$  and  $f'(x_n) = y'_n$ . This spline can be obtained by setting  $\alpha = \bar{\alpha} = -1/2$  and

$$\beta = \frac{3}{x_2 - x_1} \left[ \frac{y_2 - y_1}{x_2 - x_1} - y'_1 \right] \quad \bar{\beta} = \frac{-3}{x_n - x_{n-1}} \left[ \frac{y_n - y_{n-1}}{x_n - x_{n-1}} - y'_n \right].$$

**Programmer.** Rondall E. Jones (Sandia Laboratories).

## WEIGHTED LEAST SQUARES CUBIC SPLINE FITTING

Let  $t_1 < \dots < t_\ell$  be a sequence where  $\ell \geq 2$ ,  $\{(x_i, y_i) : i = 1, \dots, m\}$  be a set of  $m \geq 4$  points where  $t_1 \leq x_1 < \dots < x_m \leq t_\ell$ , and  $w_1, \dots, w_m$  be positive weights. Then the subroutine SPFIT is available for obtaining a cubic spline  $f(x)$  with the nodes  $t_1, \dots, t_\ell$  which minimizes  $\sum_{i=1}^m w_i (f(x_i) - y_i)^2$ . This spline is represented by

$$f(x) = z_j + a_j(x - t_j) + b_j(x - t_j)^2 + c_j(x - t_j)^3$$

for  $t_j \leq x < t_{j+1}$  ( $j = 1, \dots, \ell - 1$ ). If the nodes are selected so that  $\ell \leq m - 2$  and each interval  $(t_j, t_{j+1})$  contains a data point  $x_{\nu_j}$ , then this least squares approximation is unique.

**CALL SPFIT**( $X, Y, W, m, T, \ell, Z, A, B, C, WK, IERR$ )

$X$  is an array containing  $x_1, \dots, x_m$ ,  $Y$  an array containing  $y_1, \dots, y_m$ ,  $W$  an array containing  $w_1, \dots, w_m$ , and  $T$  an array containing  $t_1, \dots, t_\ell$ .  $Z, A, B, C$  are arrays of dimension  $\ell - 1$  or larger, and  $IERR$  is an integer variable. When SPFIT is called, if no input errors are detected the  $IERR$  is set to 0. Also, the coefficients  $z_j, a_j, b_j, c_j$  of the least squares approximating spline  $f(x)$  are computed and stored in  $Z, A, B, C$ .

$WK$  is an array of dimension  $7\ell + 18$  or larger that is a work space for the routine.

**Error Return.**  $IERR$  is set to one of the following values when an input error is detected.

$IERR = 1$  if  $\ell < 2$ .

$IERR = 2$  if  $t_1 < \dots < t_\ell$  is not satisfied.

$IERR = 3$  if  $m \geq 4$  and  $t_1 \leq x_1 \leq \dots \leq x_m \leq t_\ell$  are not satisfied.

If an error is detected, the routine immediately terminates.

**Remark.** After  $A, B, C$ , and  $Z$  have been obtained, them SCOMP may be used to evaluate the spline at any point  $x$ . SEVAL may be used if derivatives are also desired.

**Programming.** SPFIT employs the subroutines BSPP, BSL2, BSPEV, BCHFAC, and BCHSLV. SPFIT was written by A. H. Morris.

## CUBIC SPLINE EVALUATION

Given  $x_1 < \dots < x_n$ . A function  $f(x)$  is a *cubic spline having the nodes (knots)*  $x_1, \dots, x_n$  if  $f$  is a polynomial of degree  $\leq 3$  on the interval  $[x_i, x_{i+1}]$  for  $i = 1, \dots, n-1$ , and the first and second derivatives  $f'(x)$  and  $f''(x)$  exist and are continuous for all  $x$ . If  $f_i$  denotes the polynomial for the interval  $[x_i, x_{i+1}]$  then  $f_i$  has the form:

$$f_i(x) = y_i + a_i(x - x_i) + b_i(x - x_i)^2 + c_i(x - x_i)^3$$

Consequently, the function  $f$  is obtained by fitting the polynomials  $f_1, \dots, f_{n-1}$  together at the points  $x_2, \dots, x_{n-1}$ . For  $x < x_1$   $f(x) = f_1(x)$ , and for  $x > x_n$   $f(x) = f_{n-1}(x)$ . Also  $f(x_i) = y_i$  for  $i = 1, \dots, n-1$ . Hence, if  $f(x_n) = y_n$  then  $f$  interpolates the points  $(x_i, y_i)$  for  $i = 1, \dots, n$ .

A cubic spline  $f$  given by the polynomials  $f_1, \dots, f_{n-1}$  is uniquely defined by any of the following three sets of data:

- (1) the points  $(x_i, y_i)$  and coefficients  $a_i, b_i, c_i$  for  $i = 1, \dots, n-1$
- (2) the points  $(x_i, y_i)$  and first derivatives  $f'(x_i)$  for  $i = 1, \dots, n$
- (3) the points  $(x_i, y_i)$  and second derivatives  $f''(x_i)$  for  $i = 1, \dots, n$

The subroutines SCOMP, SCOMP1, SCOMP2 are available for computing the spline at any point  $x$ . SCOMP is used if data set (1) is given, SCOMP1 is used if data set (2) is given, and SCOMP2 is used if data set (3) is given.

**CALL SCOMP(X, Y, A, B, C, N, XI, YI, m, IERR)**

Let  $N = n - 1$ . Then  $N$  is the number of polynomials  $f_i$  that form the spline,  $X$  and  $Y$  are arrays containing the abscissas  $x_1, \dots, x_N$  and ordinates  $y_1, \dots, y_N$ , and  $A, B, C$  are arrays containing the coefficients  $a_i, b_i, c_i$  ( $i = 1, \dots, N$ ). It is assumed that  $N \geq 1$  and that  $x_1 < \dots < x_N$ .

Let  $\bar{x}_1, \dots, \bar{x}_m$  be the points at which the spline  $f$  is to be evaluated.  $XI$  is an array containing  $\bar{x}_1, \dots, \bar{x}_m$ ,  $YI$  an array of dimension  $m$  or larger, and  $IERR$  a variable. When SCOMP is called, if  $m < 1$  then  $IERR$  is set to 1 and the routine terminates. Otherwise, if  $m \geq 1$  then  $IERR$  is set to 0 and  $f(\bar{x}_j)$  is computed and stored in  $YI(j)$  for  $j = 1, \dots, m$ .

**Note.** SCOMP does not require  $f$  to be a spline. It is only required that  $f_i(x)$  be a cubic polynomial  $y_i + a_i(x - x_i) + b_i(x - x_i)^2 + c_i(x - x_i)^3$  and that

$$\begin{aligned} f(x) &= f_1(x) && \text{for } x < x_1 \\ f(x) &= f_i(x) && \text{for } x_i \leq x < x_{i+1} \quad (1 \leq i < N) \\ f(x) &= f_N(x) && \text{for } x \geq x_N. \end{aligned}$$

In this case, SCOMP computes the value  $f(\bar{x}_j)$  for  $j = 1, \dots, m$ .

**Programming.** Adaptation by A. H. Morris of code written by Rondall E. Jones (Sandia Laboratories).

**CALL SCOMP1**(X,Y,DY,n,XI,YI,m,IERR)

X, Y, DY are arrays containing the abscissas  $x_1, \dots, x_n$ , ordinates  $y_1, \dots, y_n$ , and first derivatives  $f'(x_1), \dots, f'(x_n)$  respectively. It is assumed that  $n \geq 2$  and  $x_1 < \dots < x_n$ .

Let  $\bar{x}_1, \dots, \bar{x}_m$  be the points at which the spline  $f$  is to be evaluated. XI is an array containing  $\bar{x}_1, \dots, \bar{x}_m$ , YI an array of dimension  $m$  or larger, and IERR a variable. When SCOMP1 is called, if  $m < 1$  then IERR is set to 1 and the routine terminates. Otherwise, if  $m \geq 1$  then IERR is set to 0 and  $f(\bar{x}_j)$  is computed and stored in YI( $j$ ) for  $j = 1, \dots, m$ .

**Programming.** Adaptation by A. H. Morris of code written by Rondall E. Jones (Sandia Laboratories).

**CALL SCOMP2**(X,Y,DDY,n,XI,YI,m,IERR)

X, Y, DDY are arrays containing the abscissas  $x_1, \dots, x_n$ , ordinates  $y_1, \dots, y_n$  and second derivatives  $f''(x_1), \dots, f''(x_n)$  respectively. It is assumed that  $n \geq 2$  and  $x_1 < \dots < x_n$ .

Let  $\bar{x}_1, \dots, \bar{x}_m$  be the points at which the spline  $f$  is to be evaluated. XI is an array containing  $\bar{x}_1, \dots, \bar{x}_m$ , YI an array of dimension  $m$  or larger, and IERR a variable. When SCOMP2 is called, if  $m < 1$  then IERR is set to 1 and the routine terminates. Otherwise, if  $m \geq 1$  then IERR is set to 0 and  $f(\bar{x}_j)$  is computed and stored in YI( $j$ ) for  $j = 1, \dots, m$ .

**Programmer.** Rondall E. Jones (Sandia Laboratories)

## CUBIC SPLINE EVALUATION AND DIFFERENTIATION

Given  $x_1 < \dots < x_n$ . A function  $f(x)$  is a *cubic spline having the nodes (knots)*  $x_1, \dots, x_n$  if  $f$  is a polynomial of degree  $\leq 3$  on the interval  $[x_i, x_{i+1}]$  for  $i = 1, \dots, n-1$ , and the first and second derivatives  $f'(x)$  and  $f''(x)$  exist and are continuous for all  $x$ . If  $f_i$  denotes the polynomial for the interval  $[x_i, x_{i+1}]$  then  $f_i$  has the form:

$$f_i(x) = y_i + a_i(x - x_i) + b_i(x - x_i)^2 + c_i(x - x_i)^3$$

Consequently, the function  $f$  is obtained by fitting the polynomials  $f_1, \dots, f_{n-1}$  together at the points  $x_2, \dots, x_{n-1}$ . For  $x < x_1$   $f(x) = f_1(x)$ , and for  $x > x_n$   $f(x) = f_{n-1}(x)$ . Also  $f(x_i) = y_i$  for  $i = 1, \dots, n-1$ . Hence, if  $f(x_n) = y_n$  then  $f$  interpolates the points  $(x_i, y_i)$  for  $i = 1, \dots, n$ .

A cubic spline  $f$  given by the polynomials  $f_1, \dots, f_{n-1}$  is uniquely defined by any of the following three sets of data:

- (1) the points  $(x_i, y_i)$  and coefficients  $a_i, b_i, c_i$  for  $i = 1, \dots, n-1$
- (2) the points  $(x_i, y_i)$  and first derivatives  $f'(x_i)$  for  $i = 1, \dots, n$
- (3) the points  $(x_i, y_i)$  and second derivatives  $f''(x_i)$  for  $i = 1, \dots, n$

The subroutines SEVAL, SEVAL1, SEVAL2 are available for computing the spline and its first and second derivatives at any point  $x$ . SEVAL is used if data set (1) is given, SEVAL1 is used if data set (2) is given, and SEVAL2 is used if data set (3) is given.

**CALL SEVAL(X,Y,A,B,C,N,XI,YI,DYI,DDYI,m,IERR)**

Let  $N = n - 1$ . Then  $N$  is the number of polynomials  $f_i$  that form the spline,  $X$  and  $Y$  are arrays containing the abscissas  $x_1, \dots, x_N$  and ordinates  $y_1, \dots, y_N$ , and  $A, B, C$  are arrays containing the coefficients  $a_i, b_i, c_i$  ( $i = 1, \dots, N$ ). It is assumed that  $N \geq 1$  and that  $x_1 < \dots < x_N$ .

Let  $\bar{x}_1, \dots, \bar{x}_m$  be the points at which the spline  $f$  and its first two derivatives are to be evaluated.  $XI$  is an array containing  $\bar{x}_1, \dots, \bar{x}_m$ ,  $YI, DYI, DDYI$  are arrays of dimension  $m$  or larger, and  $IERR$  is a variable. When SEVAL is called, if  $m < 1$  then  $IERR$  is set to 1 and the routine terminates. Otherwise, if  $m \geq 1$  then  $IERR$  is set to 0 and the values  $f(\bar{x}_j), f'(\bar{x}_j), f''(\bar{x}_j)$  are computed and stored in  $YI(j), DYI(j), DDYI(j)$  for  $j = 1, \dots, m$ .

**Note.** SEVAL does not require  $f$  to be a spline. It is only required that  $f_i(x)$  be a cubic polynomial  $y_i + a_i(x - x_i) + b_i(x - x_i)^2 + c_i(x - x_i)^3$  and that

$$\begin{aligned} f(x) &= f_1(x) & \text{for } x < x_1 \\ f(x) &= f_i(x) & \text{for } x_i \leq x < x_{i+1} \quad (1 \leq i < N) \\ f(x) &= f_N(x) & \text{for } x \geq x_N. \end{aligned}$$

In this case, SEVAL computes the values  $f(\bar{x}_j), f'(\bar{x}_j), f''(\bar{x}_j)$  for  $j = 1, \dots, m$ .

**Programming.** Adaptation by A. H. Morris of code written by Rondall E. Jones (Sandia Laboratories).

**CALL SEVAL1(X,Y,DY,n,XI,YI,DYI,DDYI,m,IERR)**

X, Y, DY are arrays containing the abscissas  $x_1, \dots, x_n$ , ordinates  $y_1, \dots, y_n$ , and first derivatives  $f'(x_1), \dots, f'(x_n)$  respectively. It is assumed that  $n \geq 2$  and  $x_1 < \dots < x_n$ .

Let  $\bar{x}_1, \dots, \bar{x}_m$  be the points at which the spline  $f$  and its first two derivatives are to be evaluated. XI is an array containing  $\bar{x}_1, \dots, \bar{x}_m$ , YI, DYI, DDYI are arrays of dimension  $m$  or larger, and IERR is a variable. When SEVAL1 is called, if  $m < 1$  then IERR is set to 1 and the routine terminates. Otherwise, if  $m \geq 1$  then IERR is set to 0 and the values  $f(\bar{x}_j), f'(\bar{x}_j), f''(\bar{x}_j)$  are computed and stored in YI(j), DYI(j), DDYI(j) for  $j = 1, \dots, m$ .

**Programming.** Adaptation by A. H. Morris of code written by Rondall E. Jones (Sandia Laboratories).

**CALL SEVAL2(X,Y,DDY,n,XI,YI,DYI,DDYI,m,IERR)**

X, Y, DDY are arrays containing the abscissas  $x_1, \dots, x_n$ , ordinates  $y_1, \dots, y_n$ , and second derivatives  $f''(x_1), \dots, f''(x_n)$  respectively. It is assumed that  $n \geq 2$  and  $x_1 < \dots < x_n$ .

Let  $\bar{x}_1, \dots, \bar{x}_m$  be the points at which the spline  $f$  and its first two derivatives are to be evaluated. XI is an array containing  $\bar{x}_1, \dots, \bar{x}_m$ , YI, DYI, DDYI are arrays of dimension  $m$  or larger, and IERR is a variable. When SEVAL2 is called, if  $m < 1$  then IERR is set to 1 and the routine terminates. Otherwise, if  $m \geq 1$  then IERR is set to 0 and the values  $f(\bar{x}_j), f'(\bar{x}_j), f''(\bar{x}_j)$  are computed and stored in YI(j), DYI(j), DDYI(j) for  $j = 1, \dots, m$ .

**Programmer.** Rondall E. Jones (Sandia Laboratories)

## INTEGRALS OF CUBIC SPLINES

Given  $x_1 < \dots < x_n$ . A function  $f(x)$  is a *cubic spline having the nodes (knots)*  $x_1, \dots, x_n$  if  $f$  is a polynomial of degree  $\leq 3$  on the interval  $[x_i, x_{i+1}]$  for  $i = 1, \dots, n-1$ , and the first and second derivatives  $f'(x)$  and  $f''(x)$  exist and are continuous for all  $x$ . If  $f_i$  denotes the polynomial for the interval  $[x_i, x_{i+1}]$  then  $f_i$  has the form:

$$f_i(x) = y_i + a_i(x - x_i) + b_i(x - x_i)^2 + c_i(x - x_i)^3$$

Consequently, the function  $f$  is obtained by fitting the polynomials  $f_1, \dots, f_{n-1}$  together at the points  $x_2, \dots, x_{n-1}$ . For  $x < x_1$   $f(x) = f_1(x)$ , and for  $x > x_n$   $f(x) = f_{n-1}(x)$ . Also  $f(x_i) = y_i$  for  $i = 1, \dots, n-1$ . Hence, if  $f(x_n) = y_n$  then  $f$  interpolates the points  $(x_i, y_i)$  for  $i = 1, \dots, n$ .

A cubic spline  $f$  given by the polynomials  $f_1, \dots, f_{n-1}$  is uniquely defined by any of the following three sets of data:

- (1) the points  $(x_i, y_i)$  and coefficients  $a_i, b_i, c_i$  for  $i = 1, \dots, n-1$
- (2) the points  $(x_i, y_i)$  and first derivatives  $f'(x_i)$  for  $i = 1, \dots, n$
- (3) the points  $(x_i, y_i)$  and second derivatives  $f''(x_i)$  for  $i = 1, \dots, n$

For any real  $\alpha$  and  $\beta$  the functions CSINT, CSINT1, CSINT2 are available for computing the integral  $\int_{\alpha}^{\beta} f(t) dt$ . CSINT is used if data set (1) is given, CSINT1 is used if data set (2) is given, and CSINT2 is used if data set (3) is given.

**CSINT**( $X, Y, A, B, C, N, \alpha, \beta$ )

Let  $N = n - 1$ . Then  $N$  is the number of polynomials  $f_i$  that form the spline,  $X$  and  $Y$  are arrays containing the abscissas  $x_1, \dots, x_N$  and ordinates  $y_1, \dots, y_N$ , and  $A, B, C$  are arrays containing the coefficients  $a_i, b_i, c_i$  ( $i = 1, \dots, N$ ). It is assumed that  $N \geq 1$  and that  $x_1 < \dots < x_N$ . Then CSINT has the value  $\int_{\alpha}^{\beta} f(t) dt$ .

**Programming.** CSINT calls the function INTRVL. CSINT was written by A. H. Morris.

**CSINT1**( $X, Y, DY, n, \alpha, \beta$ )

$X, Y, DY$  are arrays containing the abscissas  $x_1, \dots, x_n$ , ordinates  $y_1, \dots, y_n$ , and first derivatives  $f'(x_1), \dots, f'(x_n)$  respectively. It is assumed that  $n \geq 2$  and  $x_1 < \dots < x_n$ . Then CSINT1( $X, Y, DY, n, \alpha, \beta$ ) =  $\int_{\alpha}^{\beta} f(t) dt$ .

**Programming.** CSINT1 calls the function INTRVL. CSINT1 was written by A. H. Morris.

**CSINT2**( $X, Y, DDY, n, \alpha, \beta$ )

$X, Y, DDY$  are arrays containing the abscissas  $x_1, \dots, x_n$ , ordinates  $y_1, \dots, y_n$ , and second derivatives  $f''(x_1), \dots, f''(x_n)$  respectively. It is assumed that  $n \geq 2$  and  $x_1 < \dots < x_n$ . Then CSINT2( $X, Y, DDY, n, \alpha, \beta$ ) =  $\int_{\alpha}^{\beta} f(t) dt$ .

**Programming.** CSINT2 calls the function INTRVL. CSINT2 was written by A. H. Morris.



## N-DIMENSIONAL CUBIC SPLINE CLOSED CURVE FITTING

Given  $m \geq 2$  points  $x_i = (x_{i1}, \dots, x_{in})$  where  $n \geq 2$ . One procedure for fitting a closed curve to the points is to first let  $s_1 = 0, s_i = s_{i-1} + \|x_i - x_{i-1}\|$  for  $i = 2, \dots, m$ , and  $s_{m+1} = s_m + \|x_1 - x_m\|$ ,<sup>1</sup> next let  $t_i = s_i/s_{m+1}$  for  $i = 1, \dots, m+1$ , and then to find for each  $j \leq n$  the periodic cubic spline  $\gamma_j(t)$  having the knots  $t_1, \dots, t_{m+1}$  where  $\gamma_j(t_i) = x_{ij}$  for  $i \leq m$ . The mapping  $\gamma(t) = (\gamma_1(t), \dots, \gamma_n(t))$  then defines on the interval  $[0, 1]$  a closed curve which traverses the points  $(t_1, x_1), \dots, (t_m, x_m), (1, x_1)$  and satisfies  $\gamma'(0) = \gamma'(1)$  and  $\gamma''(0) = \gamma''(1)$ . For  $t < 0$  and  $t > 1$ ,  $\gamma(t)$  is defined by periodicity (having period 1). The subroutine CSLOOP is available for obtaining the derivatives  $\gamma'_j(t_i)$  which characterize this curve, the subroutine LOPCMP is available for computing the curve, and the subroutine LOPDF is available for differentiating the curve.

**CALL CSLOOP( $m, n, X, kx, T, DX, kdx, WK, IERR$ )**

$X$  is an  $m \times n$  matrix whose  $i^{\text{th}}$  row contains the point  $x_i = (x_{i1}, \dots, x_{in})$ , where  $m \geq 2$  and  $n \geq 2$ . It is assumed that the points  $x_1, \dots, x_m$  are indexed in the order that they are to be traversed by the curve  $\gamma(t)$ . It is also assumed that  $x_i \neq x_{i+1}$  for  $i = 1, \dots, m-1$  and that  $x_m \neq x_1$ .

$DX$  is a 2-dimensional array containing at least  $m$  rows and  $n$  columns. The arguments  $kx$  and  $kdx$  have the following values:

$kx$  = the number of rows in the dimension statement for  $X$  in the calling program

$kdx$  = the number of rows in the dimension statement for  $DX$  in the calling program

It is required that  $kx \geq m$  and  $kdx \geq m$ .

$IERR$  is a variable and  $T$  an array of dimension  $m$  or larger. When CSLOOP is called, if no input errors are detected then  $IERR$  is assigned the value 0 and  $t_1, \dots, t_m$  are computed and stored in  $T$ . Also, the derivatives  $\gamma'_j(t_i)$  are computed and stored in  $DX$ , where the  $i^{\text{th}}$  row of  $DX$  contains  $\gamma'(t_i) = (\gamma'_1(t_i), \dots, \gamma'_n(t_i))$ .

$WK$  is an array of dimension  $4(m-1)$  or larger that is a work space for the routine.

**Error Return.**  $IERR$  reports the following input errors:

$IERR = 1$  if  $m < 2$  or  $n < 2$

$IERR = 2$  if  $x_i = x_{i+1}$  for some  $i$ .

$IERR = 3$  if  $x_1 = x_m$ .

**Remark.** After  $T$  and  $DX$  are obtained, LOPCMP may be used to compute the curve and LOPDF may be used to differentiate the curve.

**Programming.** CSLOOP calls the subroutine CSLOP1 and function SNRM2. CSLOOP and CSLOP1 were written by A. H. Morris.

<sup>1</sup> $\|z\| = \sqrt{\sum_i z_i^2}$  for any  $z = (z_1, \dots, z_m)$ .

**CALL LOPCMP( $m, n, T, X, kx, DX, kdx, \ell, TI, Z, kz$ )**

$T$  is an array containing the knots  $t_1, \dots, t_m$ ,  $X$  an  $m \times n$  matrix whose  $i^{\text{th}}$  row contains the point  $x_i$ , and  $DX$  an  $m \times n$  matrix whose  $i^{\text{th}}$  row contains the derivative  $\gamma'(t_i)$ . The arguments  $kx$  and  $kdx$  have the following values:

$kx$  = the number of rows in the dimension statement for  $X$  in the calling program

$kdx$  = the number of rows in the dimension statement for  $DX$  in the calling program

It is assumed that  $m \geq 2$ ,  $n \geq 2$ ,  $kx \geq m$ , and  $kdx \geq m$ .

Let  $\bar{t}_1, \dots, \bar{t}_\ell$  be the points at which the curve  $\gamma$  is to be evaluated.  $TI$  is an array containing  $\bar{t}_1, \dots, \bar{t}_\ell$ , and  $Z$  a 2-dimensional array containing at least  $\ell$  rows and  $n$  columns. The argument  $kz$  is the row dimension for  $Z$  in the calling program. It is assumed that  $\ell \geq 1$  and  $kz \geq \ell$ . When LOPCMP is called,  $\gamma(\bar{t}_i) = (\gamma_1(\bar{t}_i), \dots, \gamma_n(\bar{t}_i))$  is computed and stored in the  $i^{\text{th}}$  row of  $Z$  for  $i = 1, \dots, \ell$ .

**Programmer.** A. H. Morris

**CALL LOPDF( $m, n, T, X, kx, DX, kdx, t_0, Z, DZ, DDZ$ )**

$T$  is an array containing the knots  $t_1, \dots, t_m$ ,  $X$  an  $m \times n$  matrix whose  $i^{\text{th}}$  row contains the point  $x_i$ ,  $DX$  an  $m \times n$  matrix whose  $i^{\text{th}}$  row contains the derivative  $\gamma'(t_i)$ . The arguments  $kx$  and  $kdx$  have the following values:

$kx$  = the number of rows in the dimension statement for  $X$  in the calling program

$kdx$  = the number of rows in the dimension statement for  $DX$  in the calling program

It is assumed that  $m \geq 2$ ,  $n \geq 2$ ,  $kx \geq m$ , and  $kdx \geq m$ .

$Z$ ,  $DZ$ , and  $DDZ$  are arrays of dimension  $n$ . For any real  $t_0$ ,  $\gamma(t_0) = (\gamma_1(t_0), \dots, \gamma_n(t_0))$ ,  $\gamma'(t_0) = (\gamma'_1(t_0), \dots, \gamma'_n(t_0))$ , and  $\gamma''(t_0) = (\gamma''_1(t_0), \dots, \gamma''_n(t_0))$  are computed and stored in  $Z$ ,  $DZ$ , and  $DDZ$  respectively.

**Programmer.** A. H. Morris

## SPLINE UNDER TENSION INTERPOLATION

Given real  $\sigma$  and  $x_1 < \dots < x_n$ . A function  $f(x)$  is a *spline having the tension factor  $\sigma$  and the nodes (knots)  $x_1, \dots, x_n$*  if  $f(x)$  and its first two derivatives are continuous on  $[x_1, x_n]$ , and  $f''(x) - \bar{\sigma}^2 f(x) = a_i x + b_i$  on the interval  $[x_i, x_{i+1}]$  for  $i = 1, \dots, n-1$ . Here  $\bar{\sigma} = |\sigma|(n-1)/(x_n - x_1)$  and  $a_i, b_i$  are constants. For  $x_i \leq x \leq x_{i+1}$   $f(x)$  can be represented by

$$f(x) = A_i \sinh \bar{\sigma}(x - x_i) + B_i \sinh \bar{\sigma}(x_{i+1} - x) - (a_i x + b_i)/\bar{\sigma}^2$$

when  $\sigma \neq 0$ , and by a cubic polynomial when  $\sigma = 0$ .

Assume now that  $n$  ordinates  $y_1, \dots, y_n$  are given. Then there exist an infinitude of splines  $f(x)$  having tension  $\sigma$  for which  $f(x_i) = y_i$  ( $i = 1, \dots, n$ ). However, if values  $y'_1$  and  $y'_n$  are given then only one of these splines will satisfy  $f'(x_1) = y'_1$  and  $f'(x_n) = y'_n$ . For convenience, denote this spline by  $f_\sigma$ . If  $\sigma = 0$  then it is clear that  $f_\sigma$  is the standard cubic spline. Also it can be verified that when  $\sigma \rightarrow \infty$ ,  $f_\sigma$  converges uniformly on  $[x_1, x_n]$  to the piecewise linear function  $\ell(x)$  where  $\ell(x) = y_i + m_i(x - x_i)$  for  $x_i \leq x \leq x_{i+1}$  ( $i = 1, \dots, n-1$ ). Here  $m_i = (y_{i+1} - y_i)/(x_{i+1} - x_i)$ . The following subroutine is available for obtaining the spline  $f_\sigma$ .

**CALL CURV1( $n, X, Y, SLP1, SLPN, IND, DDY, TEMP, \sigma, IERR$ )**

$X$  and  $Y$  are arrays containing the abscissas  $x_1, \dots, x_n$  and ordinates  $y_1, \dots, y_n$ . It is assumed that  $n \geq 2$  and  $x_1 < \dots < x_n$ .

$SLP1$  and  $SLPN$  are assigned the values  $y'_1$  and  $y'_n$ . The user may omit values for either or both of these arguments.  $IND$  specifies the information that is provided.

$IND = 0$  Values are supplied for  $SLP1$  and  $SLPN$ .

$IND = 1$  A value is supplied for  $SLP1$  but not for  $SLPN$ .

$IND = 2$  A value is supplied for  $SLPN$  but not for  $SLP1$ .

$IND = 3$  Values are not supplied for  $SLP1$  and  $SLPN$ .

If a value is not supplied by the user, then the routine provides a value.

$DDY$  is an array of dimension  $n$  or larger, and  $IERR$  is an integer variable. When **CURV1** is called, if no input errors are detected then  $IERR$  is assigned the value 0 and the second derivatives  $f''_\sigma(x_1), \dots, f''_\sigma(x_n)$  are computed and stored in  $DDY$ .

$TEMP$  is an array of dimension  $n$  or larger that is used for a work space.

**Error Return.**  $IERR$  reports the following input errors:

$IERR = 1$  if  $n < 2$ .

$IERR = 2$  if  $x_1 < \dots < x_n$  is not satisfied.

When either of these errors is detected, the routine immediately terminates.

**Remarks.**

- (1) After DDY is obtained then CURV2 may be used to evaluate the spline at any point  $x$ .
- (2) X, Y, n, SLP1, SLPN, IND,  $\sigma$  are not modified by CURV1.

**Programming.** CURV1 employs the subroutines CEEZ, TERMS, and SNHCSH. CURV1, CEEZ, and TERMS were written by A. K. Cline and R. J. Renka (University of Texas at Austin).

**Reference.** Cline, A. K., "Scalar and Planar Valued Curve Fitting using Splines under Tension," *Comm. ACM* 17 (1974), pp. 218-220.

## SPLINE UNDER TENSION EVALUATION

Given  $\sigma$  and  $x_1 < \dots < x_n$ . A function  $f(x)$  is a *spline having the tension factor  $\sigma$  and the nodes (knots)  $x_1, \dots, x_n$*  if  $f(x)$  and its first two derivatives are continuous on  $[x_1, x_n]$ , and  $f''(x) - \bar{\sigma}^2 f(x) = a_i x + b_i$  on the interval  $[x_i, x_{i+1}]$  for  $i = 1, \dots, n-1$ . Here  $\bar{\sigma} = |\sigma|(n-1)/(x_n - x_1)$  and  $a_i, b_i$  are constants. If  $f(x) = f_i(x)$  for  $x_i \leq x \leq x_{i+1}$  then  $f_i(x)$  can be represented by

$$f_i(x) = A_i \sinh \bar{\sigma}(x - x_i) + B_i \sinh \bar{\sigma}(x_{i+1} - x) - (a_i x + b_i)/\bar{\sigma}^2$$

when  $\sigma \neq 0$ , and by a cubic polynomial when  $\sigma = 0$ . For  $x < x_1$  we let  $f(x) = f_1(x)$ , and for  $x > x_n$  we let  $f(x) = f_{n-1}(x)$ .

Assume now that  $f(x_i) = y_i$  for  $i = 1, \dots, n$ . Then for a fixed  $\sigma$ ,  $f(x)$  is uniquely defined by the points  $(x_i, y_i)$  and the second derivatives  $f''(x_i)$  ( $i = 1, \dots, n$ ). When this data is available, the following function may be used to compute the spline at any point  $t$ .

**CURV2**( $t, n, X, Y, DDY, \sigma$ )

$X$  and  $Y$  are arrays containing the abscissas  $x_1, \dots, x_n$  and ordinates  $y_1, \dots, y_n$ , and  $DDY$  is an array containing the second derivatives  $f''(x_1), \dots, f''(x_n)$ . It is assumed that  $n \geq 2$  and  $x_1 < \dots < x_n$ .  $\text{CURV2}(t, n, X, Y, DDY, \sigma) = f(t)$  for any real  $t$ .

**Remark.** After  $DDY$  has been obtained, **CURV2** may be repeatedly called to evaluate the curve at different points so long as the tension factor  $\sigma$  remains fixed. However, if  $\sigma$  is modified then the derivative information in  $DDY$  will have to be recomputed before **CURV2** can be used with the new tension factor.

**Programming.** **CURV2** employs the function **INTRVL** and subroutine **SNHCSH**. **CURV2** was written by A. K. Cline and R. J. Renka (University of Texas at Austin). **INTRVL** was written by A. H. Morris.

**Reference.** Cline, A. K., "Scalar and Planar Valued Curve Fitting using Splines under Tension," *Comm. ACM* 17 (1974), pp. 218-220.

## DIFFERENTIATION AND INTEGRATION OF SPLINES UNDER TENSION

Let  $f(x)$  be a spline having the tension factor  $\sigma$  and the nodes  $x_1, \dots, x_n$ . Assume that  $f(x_i) = y_i$  for  $i = 1, \dots, n$ . If the second derivatives  $f''(x_1), \dots, f''(x_n)$  are known then the following functions may be used for differentiating and integrating the spline.

**CURVD**( $t, n, X, Y, DDY, \sigma$ )

$X$  and  $Y$  are arrays containing the abscissas  $x_1, \dots, x_n$  and ordinates  $y_1, \dots, y_n$ , and  $DDY$  is an array containing the second derivatives  $f''(x_1), \dots, f''(x_n)$ . It is assumed that  $n \geq 2$  and  $x_1 < \dots < x_n$ . For any real  $t$ , the derivative  $f'(t)$  is computed and assigned to be the value of **CURVD**( $t, n, X, Y, DDY, \sigma$ ).

**Programming.** **CURVD** employs the function **INTRVL** and subroutine **SNHCSH**. **CURVD** was written by A. K. Cline and R. J. Renka (University of Texas at Austin). **INTRVL** was written by A. H. Morris.

**CURVI**( $a, b, n, X, Y, DDY, \sigma$ )

$X$  and  $Y$  are arrays containing the abscissas  $x_1, \dots, x_n$  and ordinates  $y_1, \dots, y_n$ , and  $DDY$  is an array containing the second derivatives  $f''(x_1), \dots, f''(x_n)$ . It is assumed that  $n \geq 2$  and  $x_1 < \dots < x_n$ . **CURVI**( $a, b, n, X, Y, DDY, \sigma$ ) =  $\int_a^b f(t) dt$  for any real  $a$  and  $b$ .

**Note.** It is not required that  $a \leq b$ .

**Programming.** **CURVI** employs the function **INTRVL** and subroutine **SNHCSH**. **CURVI** was written by A. K. Cline and R. J. Renka (University of Texas at Austin). **INTRVL** was written by A. H. Morris.

## TWO DIMENSIONAL SPLINE UNDER TENSION CURVE FITTING

Given a sequence of points  $(x_1, y_1), \dots, (x_n, y_n)$ . One procedure for fitting a curve to the points is to let  $s_1 = 0$  and  $s_i = s_{i-1} + \sqrt{(x_i - x_{i-1})^2 + (y_i - y_{i-1})^2}$  for  $i = 2, \dots, n$ , and then to find two splines  $x(s)$  and  $y(s)$  with tension  $\sigma$  that satisfy  $x(s_i) = x_i$  and  $y(s_i) = y_i$  for  $i = 1, \dots, n$ . If  $\theta_1$  and  $\theta_n$  are the desired angles for the curve  $s \mapsto (x(s), y(s))$  at the points  $(x_1, y_1)$  and  $(x_n, y_n)$ , then the splines  $x(s)$  and  $y(s)$  can be selected so that  $x'(s_i) = \cos \theta_i$  and  $y'(s_i) = \sin \theta_i$  for  $i = 1, n$ . The curve  $s \mapsto (x(s), y(s))$  then passes through the points  $(x_i, y_i)$  and has the required slopes at the end points. The subroutine KURV1 is available for obtaining the second derivatives  $x''(s_i), y''(s_i)$  ( $i = 1, \dots, n$ ) which characterize this curve, and the subroutine KURV2 is available for computing the curve.

**CALL KURV1**( $n, X, Y, \text{SLP1}, \text{SLPN}, \text{IND}, \text{DDX}, \text{DDY}, \text{TEMP}, S, \sigma, \text{IERR}$ )

$X$  and  $Y$  are arrays containing the abscissas  $x_1, \dots, x_n$  and ordinates  $y_1, \dots, y_n$ . It is assumed that  $n \geq 2$  and that the points  $(x_i, y_i)$  are indexed in the order that they are to be traversed by the curve. It is also assumed that  $(x_i, y_i) \neq (x_{i+1}, y_{i+1})$  for  $i = 1, \dots, n-1$ .

SLP1 and SLPN are assigned the values  $\theta_1$  and  $\theta_n$ . These angles are measured counter-clockwise (in radians) from the positive x-axis. The user may omit values for SLP1 and/or SLPN. IND specifies the information that is provided.

IND = 0 Values are supplied for SLP1 and SLPN.

IND = 1 A value is supplied for SLP1 but not for SLPN.

IND = 2 A value is supplied for SLPN but not for SLP1.

IND = 3 Values are not supplied for SLP1 and SLPN.

If a value is not supplied by the user, then the routine provides a value.

$\sigma$  is the tension factor to be employed. If  $|\sigma|$  is small, say  $|\sigma| < 10^{-3}$ , then  $x(s)$  and  $y(s)$  approximate cubic splines. Otherwise, if  $|\sigma|$  is large, say  $|\sigma| > 100$ , then the resulting curve approximates the polygonal line from  $(x_1, y_1)$  to  $(x_n, y_n)$ .

IERR is an integer variable and  $S, \text{DDX}, \text{DDY}$  are arrays of dimension  $n$  or larger. When KURV1 is called, if no input errors are detected then IERR is assigned the value 0 and the values  $s_1, \dots, s_n$  are computed and stored in  $S$ . Also, the second derivatives  $x''(s_1), \dots, x''(s_n)$  and  $y''(s_1), \dots, y''(s_n)$  are computed and stored in DDX and DDY.

TEMP is an array of dimension  $n$  or larger that is used for a work space.

**Error Return.** IERR reports the following input errors:

IERR = 1 if  $n < 2$ .

IERR = 2 if  $(x_i, y_i) = (x_{i+1}, y_{i+1})$  for some  $i$ .

When either of these errors is detected, the routine immediately terminates.

**Remark.** After  $S, \text{DDX}, \text{DDY}$  are obtained, KURV2 may be used to compute the curve.

**Programming.** KURV1 employs the subroutines CEEZ, TERMS, and SNHCSH. KURV1,

CEEZ, and TERMS were written by A. K. Cline and R. J. Renka (University of Texas at Austin).

**CALL KURV2( $t, XT, YT, n, X, Y, DDX, DDY, S, \sigma$ )**

$X$  and  $Y$  are arrays containing the abscissas  $x_1, \dots, x_n$  and ordinates  $y_1, \dots, y_n$ ,  $S$  is an array containing  $s_1, \dots, s_n$ , and  $DDX$  and  $DDY$  are arrays containing the second derivatives  $x''(s_1), \dots, x''(s_n)$  and  $y''(s_1), \dots, y''(s_n)$ .

Now consider the change of variables  $t = s/s_n$ , and let  $t \mapsto (\bar{x}(t), \bar{y}(t))$  denote the curve in terms of the new parameter  $t$ .  $XT$  and  $YT$  are real variables. For any  $0 \leq t \leq 1$ , KURV2 computes the point  $(\bar{x}(t), \bar{y}(t))$  on the curve and assigns  $XT$  the value  $\bar{x}(t)$  and  $YT$  the value  $\bar{y}(t)$ .

**Remark.** After  $DDX$  and  $DDY$  have been obtained, KURV2 may be repeatedly called to evaluate the curve at different points so long as the tension factor  $\sigma$  remains fixed. However, if  $\sigma$  is modified then the derivative information in  $DDX$  and  $DDY$  will have to be recomputed before KURV2 can be used with the new tension factor.

**Programming.** KURV2 employs the function INTRVL and subroutine SNHCSH. KURV2 was written by A. K. Cline and R. J. Renka (University of Texas at Austin). INTRVL was written by A. H. Morris.



## TWO DIMENSIONAL SPLINE UNDER TENSION CLOSED CURVE FITTING

Given  $n \geq 2$  points  $(x_1, y_1), \dots, (x_n, y_n)$ . One procedure for fitting a closed curve to the points is to let  $s_1 = \sqrt{(x_1 - x_n)^2 + (y_1 - y_n)^2}$  and  $s_i = s_{i-1} + \sqrt{(x_i - x_{i-1})^2 + (y_i - y_{i-1})^2}$  for  $i = 2, \dots, n$ , and then to find periodic splines  $x(s)$  and  $y(s)$  with tension  $\sigma$  that pass through the points  $(s_1, x_1), \dots, (s_n, x_n), (s_1 + s_n, x_1)$  and  $(s_1, y_1), \dots, (s_n, y_n), (s_1 + s_n, y_1)$ . The mapping  $s \mapsto (x(s), y(s))$  then defines a closed curve that passes through the points  $(x_i, y_i)$ . The subroutine KURVP1 is available for obtaining the second derivatives  $x''(s_i), y''(s_i)$  ( $i = 1, \dots, n$ ) which characterize this curve, and the subroutine KURVP2 is available for computing the curve.

**CALL KURVP1( $n, X, Y, DDX, DDY, TEMP, S, \sigma, IERR$ )**

$X$  and  $Y$  are arrays containing the abscissas  $x_1, \dots, x_n$  and ordinates  $y_1, \dots, y_n$ . It is assumed that  $n \geq 2$  and that the points  $(x_i, y_i)$  are indexed in the order that they are to be traversed by the curve. It is also assumed that  $(x_i, y_i) \neq (x_{i+1}, y_{i+1})$  for  $i = 1, \dots, n-1$ .

$\sigma$  is the tension factor to be employed. If  $|\sigma|$  is small, say  $|\sigma| < 10^{-3}$ , then  $x(s)$  and  $y(s)$  approximate periodic cubic splines. Otherwise, if  $|\sigma|$  is large, say  $|\sigma| > 100$ , then the curve approximates the closed polygonal path that traverses the points  $(x_i, y_i)$ .

$IERR$  is an integer variable and  $S, DDX, DDY$  are arrays of dimension  $n$  or larger. When KURVP1 is called, if no input errors are detected then  $IERR$  is assigned the value 0 and the values  $s_1, \dots, s_n$  are computed and stored in  $S$ . Also, the second derivatives  $x''(s_1), \dots, x''(s_n)$  and  $y''(s_1), \dots, y''(s_n)$  are computed and stored in  $DDX$  and  $DDY$ .

$TEMP$  is an array of dimension  $2n$  or larger that is used for a work space.

**Error Return.**  $IERR$  reports the following input errors:

$IERR = 1$  if  $n < 2$ .

$IERR = 2$  if  $(x_i, y_i) = (x_{i+1}, y_{i+1})$  for some  $i$ .

When either of these errors is detected, the routine immediately terminates.

**Remark.** After  $S, DDX, DDY$  are obtained, KURVP2 may be used to compute the curve.

**Programming.** KURVP1 employs the subroutines TERMS and SNHCSH. KURVP1 and TERMS were written by A. K. Cline and R. J. Renka (University of Texas at Austin).

**CALL KURVP2( $t, XT, YT, n, X, Y, DDX, DDY, S, \sigma$ )**

$X$  and  $Y$  are arrays containing the abscissas  $x_1, \dots, x_n$  and ordinates  $y_1, \dots, y_n$ ,  $S$  is an array containing  $s_1, \dots, s_n$ , and  $DDX$  and  $DDY$  are arrays containing the second derivatives  $x''(s_1), \dots, x''(s_n)$  and  $y''(s_1), \dots, y''(s_n)$ .

Now consider the change of variables  $t = (s - s_1)/s_n$ , and let  $t \mapsto (\bar{x}(t), \bar{y}(t))$  denote the curve in terms of the new parameter  $t$ . Then  $t \mapsto (\bar{x}(t), \bar{y}(t))$  maps 0 and 1 onto the point  $(x_1, y_1)$ , and  $t \mapsto (\bar{x}(t), \bar{y}(t))$  is a periodic function (with period 1).

XT and YT are real variables. For any real  $t$ , KURVP2 computes the point  $(\bar{x}(t), \bar{y}(t))$  on the curve and assigns XT the value  $\bar{x}(t)$  and YT the value  $\bar{y}(t)$ .

**Remark.** After DDX and DDY have been obtained, KURVP2 may be repeatedly called to evaluate the curve at different points so long as the tension factor  $\sigma$  remains fixed. However, if  $\sigma$  is modified then the derivative information in DDX and DDY will have to be recomputed before KURVP2 can be used with the new tension factor.

**Programming.** KURVP2 employs the function INTRVL and subroutine SNHCSH. KURVP2 was written by A. K. Cline and R. J. Renka (University of Texas at Austin). INTRVL was written by A. H. Morris.

### THREE DIMENSIONAL SPLINE UNDER TENSION CURVE FITTING

Given  $n \geq 2$  points  $(x_1, y_1, z_1), \dots, (x_n, y_n, z_n)$ . One procedure for fitting a curve to the points is to let  $s_1 = 0$  and  $s_i = s_{i-1} + \sqrt{(x_i - x_{i-1})^2 + (y_i - y_{i-1})^2 + (z_i - z_{i-1})^2}$  for  $i = 2, \dots, n$ , and then to find splines  $x(s), y(s), z(s)$  with tension  $\sigma$  that satisfy  $x(s_i) = x_i, y(s_i) = y_i$ , and  $z(s_i) = z_i$  for  $i = 1, \dots, n$ . If  $(x'_1, y'_1, z'_1)$  and  $(x'_n, y'_n, z'_n)$  are the desired slopes for the curve  $s \mapsto (x(s), y(s), z(s))$  at the points  $(x_1, y_1, z_1)$  and  $(x_n, y_n, z_n)$ , then the splines  $x(s), y(s), z(s)$  can be selected so that  $x'(s_i) = x'_i, y'(s_i) = y'_i$ , and  $z'(s_i) = z'_i$  for  $i = 1, n$ . The curve  $s \mapsto (x(s), y(s), z(s))$  then passes through the points  $(x_i, y_i, z_i)$  and has the required slopes at the end points. The subroutine QURV1 is available for obtaining the second derivatives  $x''(s_i), y''(s_i), z''(s_i)$  ( $i = 1, \dots, n$ ) which characterize this curve, and the subroutine QURV2 is available for computing the curve.

CALL QURV1( $n, X, Y, Z, \text{SLP1X}, \text{SLP1Y}, \text{SLP1Z}, \text{SLPNX}, \text{SLPNY},$   
 $\text{SLPNZ}, \text{IND}, \text{DDX}, \text{DDY}, \text{DDZ}, \text{TEMP}, S, \sigma, \text{IERR}$ )

$X$  is an array containing  $x_1, \dots, x_n$ ,  $Y$  an array containing  $y_1, \dots, y_n$ , and  $Z$  an array containing  $z_1, \dots, z_n$ . It is assumed that  $n \geq 2$  and that the points  $(x_i, y_i, z_i)$  are indexed in the order that they are to be traversed by the curve. It is also assumed that  $(x_i, y_i, z_i) \neq (x_{i+1}, y_{i+1}, z_{i+1})$  for  $i = 1, \dots, n-1$ .

SLP1X, SLP1Y, SLP1Z and SLPNX, SLPNY, SLPNZ are assigned the values  $x'_1, y'_1, z'_1$  and  $x'_n, y'_n, z'_n$ . The user may omit values for SLP1X, SLP1Y, SLP1Z and/or SLPNX, SLPNY, SLPNZ. The argument IND specifies the information that is provided.

IND = 0 Values are supplied for SLP1X, SLP1Y, SLP1Z and SLPNX, SLPNY, SLPNZ.

IND = 1 Values are supplied for SLP1X, SLP1Y, SLP1Z but not for SLPNX, SLPNY, SLPNZ.

IND = 2 Values are supplied for SLPNX, SLPNY, SLPNZ but not for SLP1X, SLP1Y, SLP1Z.

IND = 3 No values are supplied for SLP1X, SLP1Y, SLP1Z and SLPNX, SLPNY, SLPNZ.

If a value is not supplied by the user, then the routine provides a value.

$\sigma$  is the tension factor to be employed. If  $|\sigma|$  is small, say  $|\sigma| < 10^{-3}$ , then  $x(s), y(s), z(s)$  approximate cubic splines. Otherwise, if  $|\sigma|$  is large, say  $|\sigma| > 100$ , then the resulting curve approximates the polygonal line from  $(x_1, y_1, z_1)$  to  $(x_n, y_n, z_n)$ .

IERR is an integer variable and  $S, \text{DDX}, \text{DDY}, \text{DDZ}$  are arrays of dimension  $n$  or larger. When QURV1 is called, if no input errors are detected then IERR is assigned the value 0 and the values  $s_1, \dots, s_n$  are computed and stored in  $S$ . Also, the second derivatives  $x''(s_i), y''(s_i), z''(s_i)$  ( $i = 1, \dots, n$ ) are computed and stored in  $\text{DDX}, \text{DDY}, \text{DDZ}$ .

TEMP is an array of dimension  $n$  or larger that is used for a work space.

**Error Return.** IERR reports the following input errors:

IERR = 1 if  $n < 2$ .

IERR = 2 if  $(x_i, y_i, z_i) = (x_{i+1}, y_{i+1}, z_{i+1})$  for some  $i$ .

When either of these errors is detected, the routine immediately terminates.

**Remark.** After  $S$ , DDX, DDY, DDZ are obtained, QURV2 may be used to compute the curve.

**Programming.** QURV1 employs the subroutines CEEZ, TERMS, and SNHCSH. QURV1, CEEZ, and TERMS were written by A. K. Cline and R. J. Renka (University of Texas at Austin).

**CALL QURV2**( $t, XT, YT, ZT, n, X, Y, Z, DDX, DDY, DDZ, S, \sigma$ )

$X$  is an array containing  $x_1, \dots, x_n$ ,  $Y$  an array containing  $y_1, \dots, y_n$ , and  $Z$  an array containing  $z_1, \dots, z_n$ .  $S$  is an array containing  $s_1, \dots, s_n$  and DDX, DDY, DDZ are arrays containing the second derivatives  $x''(s_i), y''(s_i), z''(s_i)$  ( $i = 1, \dots, n$ ).

Now consider the change of variables  $t = s/s_n$  and let  $t \mapsto (\bar{x}(t), \bar{y}(t), \bar{z}(t))$  denote the curve in terms of the new parameter  $t$ .  $XT, YT, ZT$  are real variables. For any  $0 \leq t \leq 1$ , QURV2 computes the point  $(\bar{x}(t), \bar{y}(t), \bar{z}(t))$  on the curve and assigns  $XT, YT, ZT$  the values  $\bar{x}(t), \bar{y}(t), \bar{z}(t)$ .

**Remark.** After DDX, DDY, DDZ have been obtained, QURV2 may be repeatedly called to evaluate the curve at different points so long as the tension factor  $\sigma$  remains fixed. However, if  $\sigma$  is modified then the derivative information in DDX, DDY, DDZ will have to be recomputed before QURV2 can be used with the new tension factor.

**Programming.** QURV2 employs the function INTRVL and subroutine SNHCSH. QURV2 was written by A. K. Cline and R. J. Renka (University of Texas at Austin).

## B-SPLINES

For  $k \geq 1$  let  $f(t) = (t-x)^{k-1}$  when  $t > x$  and  $f(t) = 0$  when  $t \leq x$ . Then for any sequence  $t_i \leq \dots \leq t_{i+k}$  where  $t_i < t_{i+k}$ , let  $B_{ik}(x) = (t_{i+k} - t_i) f[t_i, \dots, t_{i+k}]$  where  $f[t_i, \dots, t_{i+k}]$  is the  $k^{\text{th}}$  order divided difference of  $f(t)$ . The function  $B_{ik}$  is called a *B-spline of order  $k$* . For  $k = 1$  it follows that

$$B_{i1}(x) = \begin{cases} 1 & \text{if } t_i \leq x < t_{i+1} \\ 0 & \text{otherwise} \end{cases}$$

More generally, for  $k \geq 2$   $B_{ik}(x) = 0$  when  $x \notin [t_i, t_{i+k})$ . For  $t_i \leq x < t_{i+k}$

$$B_{ik}(x) = \left( \frac{t_{i+k} - x}{t_{i+k} - t_i} \right)^{k-1} \quad \text{when } t_i = \dots = t_{i+k-1} \text{ and}$$

$$B_{ik}(x) = \left( \frac{x - t_i}{t_{i+k} - t_i} \right)^{k-1} \quad \text{when } t_{i+1} = \dots = t_{i+k}.$$

Otherwise, if no point appears more than  $k-1$  times in the sequence  $\{t_i, \dots, t_{i+k}\}$  then

$$B_{ik}(x) = \frac{x - t_i}{t_{i+k-1} - t_i} B_{i,k-1}(x) + \frac{t_{i+k} - x}{t_{i+k} - t_{i+1}} B_{i+1,k-1}(x).$$

From these relations it follows that  $B_{ik}(x) > 0$  when  $t_i < x < t_{i+k}$ . Now let  $\xi_1 < \dots < \xi_r$  be the distinct points in  $\{t_i, \dots, t_{i+k}\}$ , where  $\xi_j$  appears  $m_j$  times for  $j = 1, \dots, r$ . Then it can be verified that  $B_{ik}$  is a polynomial of order  $\leq k$  (degree  $\leq k-1$ ) on each interval  $[\xi_j, \xi_{j+1})$  ( $j = 1, \dots, r-1$ ), and that  $B_{ik}$  is of class  $C^{k-m_j-1}$  at  $\xi_j$  for  $j = 1, \dots, r$ .

We note in passing that if  $t_i$  ( $i = 0, \pm 1, \pm 2, \dots$ ) is a sequence where  $t_i \leq t_{i+1}$  and  $t_i < t_{i+k}$  for each  $i$ , then for any  $x \in [t_j, t_{j+1})$ ,  $B_{ik}(x) \neq 0$  only when  $i = j - k + 1, \dots, j$ . Moreover, it can be verified that  $\sum_i B_{ik}(x) = 1$ .

Now let  $\xi_1 < \dots < \xi_{\ell+1}$  be a sequence of points, which we shall call *knots* or *break points*. If  $\ell \geq 2$  then  $\xi_2, \dots, \xi_{\ell}$  will be called the *interior knots*. For each interior knot  $\xi_j$  let there be associated an integer  $m_j \geq 1$ , called the *multiplicity* of the knot. Then for any  $k \geq \max \{m_2, \dots, m_{\ell}\}$  let  $t_1 \leq \dots \leq t_{n+k}$  ( $n = k + m_2 + \dots + m_{\ell}$ ) be any sequence where

- (1)  $t_1 \leq \dots \leq t_k = \xi_1$ ,
- (2)  $t_{k+1}, \dots, t_n$  are the interior knots, where each interior knot  $\xi_j$  appears exactly  $m_j$  times, and
- (3)  $\xi_{\ell+1} = t_{n+1} \leq \dots \leq t_{n+k}$ .

Otherwise, if  $\ell = 1$  then for  $k \geq 1$  let  $t_1 \leq \dots \leq t_{n+k}$  ( $n = k$ ) be any sequence where  $t_1 \leq \dots \leq t_k = \xi_1$  and  $\xi_2 = t_{n+1} \leq \dots \leq t_{n+k}$ . Then for  $\ell \geq 1$ , we note that  $B_{1k}, \dots, B_{nk}$  are the only B-splines of order  $k$  which need not be 0 on the interval  $[t_k, t_{n+1})$ . Let  $\mathcal{P}_k[t_k, \dots, t_{n+1}]$  denote the collection of all piecewise polynomials  $p(x)$  defined on the interval  $[t_k, t_{n+1})$  where  $p(x)$  is a polynomial of order  $\leq k$  (degree  $\leq k-1$ ) on  $[\xi_j, \xi_{j+1})$

for  $j = 1, \dots, \ell$ , and  $p(x)$  is of class  $C^{k-m_j-1}$  at each interior knot  $\xi_j$  where  $m_j < k$ . Then by the above remarks  $B_{1k}, \dots, B_{nk}$  are in  $\mathcal{P}_k[t_k, \dots, t_{n+1}]$ . Also it can be verified that  $B_{1k}, \dots, B_{nk}$  form a basis for the vector space  $\mathcal{P}_k[t_k, \dots, t_{n+1}]$ . Thus any piecewise polynomial  $p(x)$  in  $\mathcal{P}_k[t_k, \dots, t_{n+1}]$  can be represented uniquely in the form  $p(x) = \sum_{i=1}^n a_i B_{ik}(x)$  for  $t_k \leq x < t_{n+1}$ . This representation is called a *B-spline representation* for  $p(x)$ .

## PIECEWISE POLYNOMIAL INTERPOLATION

For  $n \geq k \geq 1$  let  $t_1 \leq \dots \leq t_{n+k}$  be a sequence where  $t_i < t_{i+k}$  for  $i = 1, \dots, n$ . Further assume that  $t_k < t_{k+1}$  and  $t_n < t_{n+1}$ , and consider a set of  $n$  points  $\{(x_i, y_i) : i = 1, \dots, n\}$  where  $t_k \leq x_1 < \dots < x_n \leq t_{n+1}$ . Then we wish to find a piecewise polynomial  $f = \sum_{i=1}^n a_i B_{ik}$  defined on the interval  $[t_k, t_{n+1})$  which satisfies  $f(x_i) = y_i$  for  $i = 1, \dots, n$ . [If  $x_n = t_{n+1}$ , then by  $f(x_n) = y_n$  we mean  $f(x_n -) = y_n$ .] This problem has a unique solution when  $x_1 < t_{k+1}$ ,  $t_i < x_i < t_{i+k}$  for  $1 < i < n$ , and  $x_n > t_n$ . The following subroutine is available for obtaining the coefficients  $a_1, \dots, a_n$  of the interpolating piecewise polynomial.

**CALL BSTRP(X, Y, T, n, k, A, WK, IFLAG)**

$X$  is an array containing  $x_1, \dots, x_n$ ,  $Y$  an array containing  $y_1, \dots, y_n$ , and  $T$  an array containing  $t_1, \dots, t_{n+k}$ .  $A$  is an array of dimension  $n$  or larger, and  $IFLAG$  an integer variable. On an initial call to the routine the user may assign  $IFLAG$  any nonzero value. In this case, if no errors are detected then  $IFLAG$  is reset by the routine to 0 and the B-spline coefficients  $a_1, \dots, a_n$  are computed and stored in  $A$ . The routine may be recalled with  $IFLAG = 0$  on input when only  $Y$  is modified. In this case, no error checking is performed and  $IFLAG = 0$  on output. Also the B-spline coefficients  $a_1, \dots, a_n$  of the new interpolating piecewise polynomial are computed and stored in  $A$ .

$WK$  is an array of dimension  $(2k - 1)n$  or larger that is used for temporary storage by the routine. When **BSTRP** terminates,  $WK$  contains information needed for subsequent calls to the routine.

**Error Return.**  $IFLAG$  is assigned the value 1 if a violation of any of the conditions

$$\begin{aligned} x_1 &< \dots < x_n \\ t_k &\leq x_1 < t_{k+1} \\ t_i &< x_i < t_{i+k} \text{ for } 1 < i < n \\ t_n &< x_n \leq t_{n+1} \end{aligned}$$

is detected. When an error is detected, the routine immediately terminates.

**Example.** Given  $n > 4$  data points  $(x_i, y_i)$ , then for  $k = 4$  one may set  $t_1 = \dots = t_k = x_1$ ,  $t_{k+i} = x_{i+2}$  for  $i = 1, \dots, n - k$ , and  $x_n = t_{n+1} = \dots = t_{n+k}$ . Then  $x_3, \dots, x_{n-2}$  are the interior knots for the interpolating piecewise polynomial  $f$ . Here we have cubic spline interpolation where the data points  $x_2$  and  $x_{n-1}$  are not knots for  $f$ .

**Selection of  $t_i \leq \dots \leq t_{n+k}$  given the data  $(x_i, y_i)$ .** It is recommended that one set  $t_1 = \dots = t_k$  and  $t_{n+1} = \dots = t_{n+k}$ . For  $k \geq 2$  it is frequently convenient to select  $n - k$  points in  $x_2, \dots, x_{n-1}$  to be the interior knots for  $f$ . (This was done in the above example.) If  $k > 2$  then an alternative approach, which often gives excellent results, is to set  $t_{k+i} = (x_{i+1} + \dots + x_{i+k-1}) / (k - 1)$  for  $i = 1, \dots, n - k$ .

**Remark.** After the B-spline representation  $\sum a_i B_{ik}$  is obtained, then the subroutine BSPP can be used to obtain the Taylor series representation. The Taylor series representation is what is normally used for evaluating piecewise polynomials.

**Programming.** BSTRP calls the subroutines BSPEV, BANFAC, and BANSLV. BSTRP is a modified version by A. H. Morris of the subroutine SPLINT. The routines SPLINT, BANFAC, and BANSLV were written by Carl de Boor (University of Wisconsin).

**Reference.** de Boor, Carl, *A Practical Guide to Splines*, Springer-Verlag, 1978.



## CONVERSION OF PIECEWISE POLYNOMIALS FROM B-SPLINE TO TAYLOR SERIES FORM

For  $n \geq k \geq 1$  let  $t_1 \leq \dots \leq t_{n+k}$  be a sequence where  $t_i < t_{i+k}$  for  $i = 1, \dots, n$ . Further assume that  $t_k < t_{n+1}$  and let  $f(x) = \sum_{i=1}^n a_i B_{i,k}(x)$  for  $t_k \leq x < t_{n+1}$ . If  $\xi_1 < \dots < \xi_{\ell+1}$  are the distinct points in the sequence  $\{t_k, \dots, t_{n+1}\}$  then the piecewise polynomial  $f$  can be represented in the form  $f(x) = \sum_{i=1}^k c_{ij}(x - \xi_j)^{i-1}$  for  $\xi_j \leq x < \xi_{j+1}$  ( $j = 1, \dots, \ell$ ). The following subroutine is available for obtaining the coefficients  $c_{ij}$  of this representation.

**CALL BSPP(T, A, n, k, BREAK, C, L, WK)**

$T$  is an array containing  $t_1, \dots, t_{n+k}$  and  $A$  an array containing  $a_1, \dots, a_n$ . **BREAK** is an array of dimension  $\ell + 1$  or larger,  $C$  a 2-dimensional array of dimension  $k \times \ell$ , and  $L$  a variable. When **BSPP** is called then  $L$  is assigned the value  $\ell$  (which is computed by the routine), the break points  $\xi_1 < \dots < \xi_{\ell+1}$  are found and stored in **BREAK**, and the coefficients  $c_{ij}$  are computed and stored in  $C$ . The  $j^{\text{th}}$  column of the matrix  $C$  then contains the coefficients of the  $j^{\text{th}}$  polynomial forming  $f$  ( $j = 1, \dots, \ell$ ).

**WK** is an array of dimension  $k(k+1)$  or larger that is used for work space by the routine.

### Remarks.

- (1) Since  $\ell \leq n - k + 1$ , **BREAK** may be declared to be of dimension  $n - k + 2$  and  $C$  to be of dimension  $k \times (n - k + 1)$ .
- (2) After  $C$  is obtained, then **PPVAL** may be used to evaluate  $f$  at any point  $x$ .

**Programming.** **BSPP** is a reformulation by A. H. Morris of the subroutine **BSPLPP**, written by Carl de Boor (University of Wisconsin).

**Reference.** de Boor, Carl, *A Practical Guide to Splines*, Springer-Verlag, 1978.

## PIECEWISE POLYNOMIAL EVALUATION

Given  $x_1 < \dots < x_{\ell+1}$ . If  $f$  is a piecewise polynomial where  $f(x) = \sum_{i=1}^k c_{ij}(x - x_j)^{i-1}$  for  $x_j \leq x < x_{j+1}$  ( $j = 1, \dots, \ell$ ), then the following subroutine is available for computing  $f$  at any point  $x$ .

**CALL PPVAL**( $X, C, k, \ell, XI, YI, m$ )

$X$  is an array containing the knots  $x_1, \dots, x_\ell$  and  $C$  a  $k \times \ell$  matrix containing the coefficients  $c_{ij}$ . It is assumed that  $k \geq 1$  and  $\ell \geq 1$ . Let  $\bar{x}_1, \dots, \bar{x}_m$  be the points at which  $f$  is to be evaluated.  $XI$  is an array containing  $\bar{x}_1, \dots, \bar{x}_m$  and  $YI$  an array of dimension  $m$  or larger. When PPVAL is called,  $f(\bar{x}_j)$  is computed and stored in  $YI(j)$  for  $j = 1, \dots, m$ .

### Remarks.

- (1)  $X$  need not contain the knot  $x_{\ell+1}$ .
- (2) It is not required that  $f$  be continuous at an interior knot  $x_i$ . If  $x_i$  appears in  $XI$  then  $f(x_i+)$  is computed.
- (3) It is not required that the output points  $\bar{x}_j$  in  $XI$  be in the interval  $[x_1, x_{\ell+1})$ . If  $\bar{x}_j < x_1$  then  $\sum_i c_{i1}(x - x_1)^{i-1}$  is evaluated at  $\bar{x}_j$ . Otherwise, if  $\bar{x}_j \geq x_{\ell+1}$  then  $\sum_i c_{i\ell}(x - x_\ell)^{i-1}$  is evaluated at  $\bar{x}_j$ .

**Programming.** PPVAL is an adaptation by A. H. Morris of code written by Rondall E. Jones (Sandia Laboratories).

## WEIGHTED LEAST SQUARES PIECEWISE POLYNOMIAL FITTING

For  $n \geq k \geq 1$  let  $t_1 \leq \dots \leq t_{n+k}$  be a sequence where  $t_i < t_{i+k}$  for  $i = 1, \dots, n$ . Further assume that  $t_k < t_{k+1}$  and  $t_n < t_{n+1}$ , and consider a set of points  $\{(x_i, y_i) : i = 1, \dots, m\}$  where  $t_k \leq x_1 \leq \dots \leq x_m \leq t_{n+1}$ . Let  $w_i > 0$  ( $i = 1, \dots, m$ ) be weights. Then the subroutine BSL2 is available for finding a piecewise polynomial  $f = \sum_{i=1}^n a_i B_{ik}$  defined on the interval  $[t_k, t_{n+1})$  which minimizes  $\sum_{i=1}^m w_i (f(x_i) - y_i)^2$ .<sup>1</sup>

**CALL BSL2**( $T, n, k, X, Y, WGT, m, A, WK, Q, IERR$ )

$T$  is an array containing  $t_1, \dots, t_{n+k}$ ,  $X$  an array containing  $x_1, \dots, x_m$ ,  $Y$  an array containing  $y_1, \dots, y_m$ , and  $WGT$  an array containing  $w_1, \dots, w_m$ .  $A$  is an array of dimension  $n$  or larger, and  $IERR$  an integer variable. When BSL2 is called, if no input errors are detected then  $IERR$  is set to 0 and the B-spline coefficients  $a_1, \dots, a_n$  of the least squares approximation  $f$  are computed and stored in  $A$ .

$WK$  is an array of dimension  $n$  or larger, and  $Q$  an array of dimension  $kn$  or larger.  $WK$  and  $Q$  are work spaces for the routine.

**Error Return.**  $IERR$  is assigned the value 1 if any of the conditions

$$\begin{aligned} n &\geq k \geq 1 \\ t_n &< t_{n+1} \\ t_k &\leq x_1 \leq \dots \leq x_m \leq t_{n+1} \end{aligned}$$

is violated. When an error is detected, the routine immediately terminates.

**Selection of  $t_1 \leq \dots \leq t_{n+k}$  given the data  $(x_i, y_i)$ .** It is recommended that the knots  $t_\nu$  be selected so that there are data points  $x_{i_1} < \dots < x_{i_n}$  satisfying

$$\begin{aligned} t_k &\leq x_{i_1} < t_{k+1} \\ t_\nu &< x_{i_\nu} < t_{\nu+k} \text{ for } \nu = 2, \dots, n. \end{aligned}$$

If these conditions are satisfied then the least squares approximation is unique.

**Remark.** After the B-spline representation  $\sum_i a_i B_{ik}$  of the least squares approximation is obtained, then the subroutine BSPP can be used to obtain the Taylor series representation. The Taylor series representation is normally used for evaluating piecewise polynomials.

**Programming.** BSL2 calls the subroutines BSPEV, BCHFAC, and BCHSLV. BSL2 is a modified version by A. H. Morris of the subroutine L2APPR. L2APPR, BCHFAC, and BCHSLV were written by Carl de Boor (University of Wisconsin).

**Reference.** de Boor, Carl, *A Practical Guide to Splines*, Springer-Verlag, 1978.

<sup>1</sup>If  $x_i = t_{n+1}$  then by  $f(x_i)$  we mean  $f(x_i^-)$ .

## BI-SPLINES UNDER TENSION

Given a tension factor  $\sigma$  and a set  $X = \{x_1, \dots, x_m\}$  where  $x_1 < \dots < x_m$ . Let  $S_\sigma(X)$  denote the collection of splines having tension  $\sigma$  and the knots  $x_1, \dots, x_m$ . Then  $S_\sigma(X)$  is a vector space. Also each spline  $f(x)$  in  $S_\sigma(X)$  is uniquely characterized by the values  $f(x_1), \dots, f(x_m)$  and the slopes  $f'(x_1)$  and  $f'(x_m)$ . Let  $\psi_i(x)$  ( $i = 1, \dots, m$ ) denote the spline in  $S_\sigma(X)$  satisfying

$$\begin{aligned}\psi_i(x_i) &= 1 \\ \psi_i(x_k) &= 0 \quad \text{for } k \neq i \\ \psi'_i(x_1) &= \psi'_i(x_m) = 0\end{aligned}$$

and let  $\psi_{m+1}, \psi_{m+2}$  be the splines satisfying

$$\begin{aligned}\psi_{m+1}(x_i) &= 0 \quad \psi_{m+2}(x_i) = 0 \quad (i = 1, \dots, m) \\ \psi'_{m+1}(x_1) &= 1 \quad \psi'_{m+2}(x_1) = 0 \\ \psi'_{m+1}(x_m) &= 0 \quad \psi'_{m+2}(x_m) = 1.\end{aligned}$$

Then  $\{\psi_1, \dots, \psi_{m+2}\}$  is a basis for  $S_\sigma(X)$  and  $f = \sum_{i=1}^m f(x_i)\psi_i + f'(x_1)\psi_{m+1} + f'(x_m)\psi_{m+2}$  for each spline  $f(x)$  in  $S_\sigma(X)$ .

Let  $Y = \{y_1, \dots, y_n\}$  where  $y_1 < \dots < y_n$  and let  $\{\bar{\psi}_1, \dots, \bar{\psi}_{n+2}\}$  be the corresponding basis for  $S_\sigma(Y)$ . Then we note that there exists a unique surface

$$F(x, y) = \sum_{i=1}^{m+2} \sum_{j=1}^{n+2} a_{ij} \bar{\psi}_i(x) \bar{\psi}_j(y)$$

in the tensor product space  $S_\sigma(X) \otimes S_\sigma(Y)$  which satisfies the conditions

$$(*) \quad \left. \begin{aligned} F(x_i, y_j) &= f(x_i, y_j) & i = 1, \dots, m \quad j = 1, \dots, n \\ D_2 F(x_i, y_1) &= D_2 f(x_i, y_1) \\ D_2 F(x_i, y_n) &= D_2 f(x_i, y_n) \end{aligned} \right\} \quad i = 1, \dots, m$$

$$\left. \begin{aligned} D_1 F(x_1, y_j) &= D_1 f(x_1, y_j) \\ D_1 F(x_m, y_j) &= D_1 f(x_m, y_j) \end{aligned} \right\} \quad j = 1, \dots, n$$

$$D_1 D_2 F(x_k, y_\ell) = D_1 D_2 f(x_k, y_\ell) \quad k = 1, m \quad \ell = 1, n$$

for given data  $f(x_i, y_j), D_2 f(x_i, y_1), \dots, D_1 D_2 f(x_k, y_\ell)$ .<sup>1</sup> Such a surface is called a **bi-spline with tension  $\sigma$** . It is easily checked that the bi-spline  $F(x, y)$  has the following properties:

- (1)  $F(x, y)$  is a  $C^2$  mapping on  $[x_1, x_m] \times [y_1, y_n]$ .
- (2) The partial derivatives  $D_1^2 D_2^2 F(x, y)$  and  $D_2^2 D_1^2 F(x, y)$  exist and are continuous, and  $D_1^2 D_2^2 F(x, y) = D_2^2 D_1^2 F(x, y)$ .
- (3) For each fixed  $y$  the mapping  $x \mapsto F(x, y)$  is a spline in  $S_\sigma(X)$ , and for each fixed  $x$  the mapping  $y \mapsto F(x, y)$  is a spline in  $S_\sigma(Y)$ .

<sup>1</sup>  $D_1 F$  and  $D_2 F$  denote the partial derivatives of  $F$ .

For a given tension factor  $\sigma$ , let  $F_\sigma$  denote the unique bi-spline in  $S_\sigma(X) \otimes S_\sigma(Y)$  which satisfies conditions (\*). If  $\sigma = 0$  then  $F_\sigma$  is the standard bicubic spline. Also it can be verified that when  $\sigma \rightarrow \infty$ ,  $F_\sigma$  converges uniformly on  $[x_1, x_m] \times [y_1, y_n]$  to the piecewise bilinear function  $\ell(x, y)$  where

$$\begin{aligned} \ell(x, y) = & f(x_i, y_j) \frac{x_{i+1} - x}{h_i} \frac{y_{j+1} - y}{k_j} + f(x_i, y_{j+1}) \frac{x_{i+1} - x}{h_i} \frac{y - y_j}{k_j} \\ & + f(x_{i+1}, y_j) \frac{x - x_i}{h_i} \frac{y_{j+1} - y}{k_j} + f(x_{i+1}, y_{j+1}) \frac{x - x_i}{h_i} \frac{y - y_j}{k_j} \end{aligned}$$

for  $x_i \leq x \leq x_{i+1}$  and  $y_j \leq y \leq y_{j+1}$ . Here  $h_i = x_{i+1} - x_i$  and  $k_j = y_{j+1} - y_j$ .

## BI-SPLINE UNDER TENSION SURFACE INTERPOLATION

Given  $x_1 < \dots < x_m$  and  $y_1 < \dots < y_n$ . Also assume that we are given the values  $z_{ij}$  ( $i = 1, \dots, m; j = 1, \dots, n$ ) and a tension factor  $\sigma$ . Then the subroutine SURF is available for finding a bi-spline  $F(x, y)$  with tension  $\sigma$  that satisfies  $F(x_i, y_j) = z_{ij}$  for each  $i, j$ . Boundary conditions can be imposed on the surface  $F(x, y)$  if desired.

**CALL SURF**( $m, n, X, Y, Z, kz, OPT, DDZ, WK, \sigma, IERR$ )

$X$  is an array containing  $x_1, \dots, x_m$ ,  $Y$  an array containing  $y_1, \dots, y_n$ , and  $Z$  the  $m \times n$  matrix  $(z_{ij})$ . The argument  $kz$  is the number of rows in the dimension statement for  $Z$  in the calling program. It is assumed that  $m \geq 2, n \geq 2$ , and  $kz \geq m$ .

$OPT$  is an array, called the *option vector*, which permits the user to specify any boundary conditions that are to be imposed on the surface. If no boundary conditions are to be specified then  $OPT$  may be declared to be of dimension 1 and  $OPT(1)$  must be assigned the value 0. The details concerning the specification of boundary conditions in  $OPT$  are given below.

$DDZ$  is a 3-dimensional array of dimension  $m \times n \times 3$  and  $IERR$  is a variable. When SURF is called, if no input errors are detected then  $IERR$  is assigned the value 0 and the partial derivatives  $D_1^2 F(x_i, y_j), D_2^2 F(x_i, y_j), D_1 D_2 F(x_i, y_j)$  ( $i = 1, \dots, m; j = 1, \dots, n$ ) are computed and stored in  $DDZ$ .  $DDZ(i, j, 1) = D_2^2 F(x_i, y_j)$ ,  $DDZ(i, j, 2) = D_1^2 F(x_i, y_j)$ , and  $DDZ(i, j, 3) = D_1 D_2 F(x_i, y_j)$  for each  $i, j$ .

$WK$  is an array of dimension  $m + 2n$  or larger that is used for a work space.

**Error Return.**  $IERR$  reports the following input errors:

$IERR = 1$  if  $m < 2$  or  $n < 2$ .

$IERR = 2$  if  $x_1 < \dots < x_m$  or  $y_1 < \dots < y_n$  is not satisfied.

$IERR = 3$  if  $OPT$  contains an error.

When an error is detected, the routine immediately terminates.

**Remark.** After  $DDZ$  is obtained then SURF2 and NSURF2 may be used to evaluate the bi-spline  $F(x, y)$ .

**The option vector  $OPT$ .** If no boundary conditions are to be imposed then  $OPT$  may be declared to be of dimension 1 and  $OPT(1)$  must have the value 0. Otherwise,  $OPT$  is an array containing the information  $key_1, data_1, key_2, data_2, \dots, key_s, data_s, 0$ . The last entry in  $OPT$  is the value 0. Each group of data  $key_i, data_i$  ( $i = 1, \dots, s$ ) is called an *option*. Each  $key_i$  is an integer and  $data_i$  is a list of partial derivative values that the surface  $F(x, y)$  is required to satisfy. The following options are available:

key = 1 The values  $D_1 F(x_1, y_j)$  ( $j = 1, \dots, n$ ) must be satisfied.

key = 2 The values  $D_1 F(x_m, y_j)$  ( $j = 1, \dots, n$ ) must be satisfied.

key = 3 The values  $D_2 F(x_i, y_1)$  ( $i = 1, \dots, m$ ) must be satisfied.

key = 4 The values  $D_2 F(x_i, y_n)$  ( $i = 1, \dots, m$ ) must be satisfied.

- key = 5 The value  $D_1 D_2 F(x_1, y_1)$  must be satisfied.
- key = 6 The value  $D_1 D_2 F(x_m, y_1)$  must be satisfied.
- key = 7 The value  $D_1 D_2 F(x_1, y_n)$  must be satisfied.
- key = 8 The value  $D_1 D_2 F(x_m, y_n)$  must be satisfied.

The order of the options in OPT is arbitrary. If an unrecognized key is used then the error indicator IERR is assigned the value 3 and the routine terminates.

**Example.** Assume that we have an array DY1 containing values  $D_2 F(x_i, y_1)$  ( $i = 1, \dots, m$ ) which are to be satisfied, and that we also want  $D_1 D_2 F(x_m, y_n) = -1.3$  to be satisfied. Then OPT must be of dimension  $\geq m + 4$  and OPT can be defined as follows:

```

OPT(1) = 3.0          (First option)
DO 10 I = 1,M
10 OPT (I + 1) = DY1(I)
OPT (M + 2) = 8.0     (Second option)
OPT (M + 3) = -1.3
OPT (M + 4) = 0.0     (Terminates the option vector)
```

**Background.** The evaluation of  $D_1^2 F(x_i, y_j)$ ,  $D_2^2 F(x_i, y_j)$ , and  $D_1^2 D_2^2 F(x_i, y_j)$  reduce to the evaluation of second derivatives of splines. Specifically, for each  $i \leq m$   $D_2^2 F(x_i, y_1), \dots, D_2^2 F(x_i, y_n)$  are the second derivatives that characterize the spline  $y \mapsto F(x_i, y)$ , and for each  $j \leq n$   $D_1^2 F(x_1, y_j), \dots, D_1^2 F(x_m, y_j)$  are the second derivatives that characterize the spline  $x \mapsto F(x, y_j)$ . Also  $D_1 D_2^2 F(x_1, y_j)$  and  $D_1 D_2^2 F(x_m, y_j)$  ( $j = 1, \dots, n$ ) are the second derivatives that characterize the splines  $y \mapsto D_1 F(x_1, y)$  and  $y \mapsto D_1 F(x_m, y)$ . For each  $j \leq n$ , after one obtains the values  $D_2^2 F(x_i, y_j)$  through which the spline  $x \mapsto D_2^2 F(x, y_j)$  will pass and the end slopes  $D_1 D_2^2 F(x_1, y_j)$  and  $D_1 D_2^2 F(x_m, y_j)$  which this spline must have, then the second derivatives that characterize this spline can be computed.  $D_1^2 D_2^2 F(x_1, y_j), \dots, D_1^2 D_2^2 F(x_m, y_j)$  are the second derivatives that characterize  $x \mapsto D_2^2 F(x, y_j)$ .

**Programming.** SURF employs the subroutines CEEZ, TERMS, and SNHCSH. SURF was written by A. K. Cline and R. J. Renka (University of Texas at Austin), and modified by A. H. Morris.

## BI-SPLINE UNDER TENSION EVALUATION

Given  $x_1 < \dots < x_m$  and  $y_1 < \dots < y_n$ , and let  $F(x, y)$  be a bi-spline with tension  $\sigma$ . If the partial derivatives  $D_1^2 F(x_i, y_j)$ ,  $D_2^2 F(x_i, y_j)$ ,  $D_1^2 D_2^2 F(x_i, y_j)$  are known for  $i = 1, \dots, m$  and  $j = 1, \dots, n$ , then the function SURF2 may be used for evaluating  $F(x, y)$  at a single point, and the subroutine NSURF2 may be used for evaluating  $F(x, y)$  on a grid of points.

**SURF2**( $s, t, m, n, X, Y, Z, kz, DDZ, \sigma$ )

$X$  is an array containing  $x_1, \dots, x_m$ ,  $Y$  an array containing  $y_1, \dots, y_n$ , and  $Z$  an  $m \times n$  matrix containing the values  $F(x_i, y_j)$ . The argument  $kz$  is the number of rows in the dimension statement for  $Z$  in the calling program. It is assumed that  $m \geq 2$ ,  $n \geq 2$ , and  $kz \geq m$ .

$DDZ$  is a 3-dimensional array of dimension  $m \times n \times 3$  containing the partial derivatives where

$$DDZ(i, j, 1) = D_2^2 F(x_i, y_j)$$

$$DDZ(i, j, 2) = D_1^2 F(x_i, y_j)$$

$$DDZ(i, j, 3) = D_1^2 D_2^2 F(x_i, y_j)$$

for each  $i, j$ .  $SURF2(s, t, m, n, X, Y, Z, kz, DDZ, \sigma) = F(s, t)$  for any point  $(s, t)$ .

**Remark.** After  $DDZ$  has been obtained,  $SURF2$  may be repeatedly called to evaluate the surface at different points so long as the tension factor  $\sigma$  remains fixed. However, if  $\sigma$  is modified then the derivative information in  $DDZ$  will have to be recomputed before  $SURF2$  can be used with the new tension factor.

**Programming.**  $SURF2$  employs the function INTRVL and subroutine SNHCSH.  $SURF2$  was written by A. K. Cline and R. J. Renka (University of Texas at Austin).

**CALL NSURF2**( $s_{\min}, s_{\max}, m_s, t_{\min}, t_{\max}, n_t, W, kw, m, n,$   
 $X, Y, Z, kz, DDZ, WORK, \sigma$ )

The arguments  $s_{\min}$  and  $s_{\max}$  are the lower and upper limits of the  $x$ -coordinates of the grid on which  $F(x, y)$  is to be evaluated, and the arguments  $t_{\min}$  and  $t_{\max}$  are the lower and upper limits of the  $y$ -coordinates. The purpose of the routine is to evaluate the surface at the points  $(s_i, t_j)$  where

$$s_i = s_{\min} + (i - 1) \frac{s_{\max} - s_{\min}}{m_s - 1}$$

$$t_j = t_{\min} + (j - 1) \frac{t_{\max} - t_{\min}}{n_t - 1}$$

for  $i = 1, \dots, m_s$  and  $j = 1, \dots, n_t$ . It is assumed that  $m_s \geq 1$  and  $n_t \geq 1$ .

$W$  is a 2-dimensional array of dimension  $kw \times n_t$  where  $kw \geq m_s$ . When  $NSURF2$  is called  $W(i, j)$  is assigned the value  $F(s_i, t_j)$  for  $i = 1, \dots, m_s$  and  $j = 1, \dots, n_t$ .



The arguments  $m, n, X, Y, Z, kz, DDZ, \sigma$  are the same as in SURF2. WORK is an array of dimension  $4m$ , or larger that is used for a work space.

**Programming.** NSURF2 employs the subroutine SNHCSH. NSURF2 was written by A. K. Cline (University of Texas at Austin).

## SURFACE INTERPOLATION FOR ARBITRARILY POSITIONED DATA POINTS

Let  $\{(x_i, y_i, z_i) : i = 1, \dots, n\}$  be a set of 4 or more points which are not collinear. If  $(x_i, y_i) \neq (x_j, y_j)$  for  $i \neq j$  then the problem is to find a smooth mapping  $z = F(x, y)$  for which  $z_i = F(x_i, y_i)$  for  $i = 1, \dots, n$ . The desired degree of smoothness might vary, but it is almost always required that  $F(x, y)$  be at least continuously differentiable.

A procedure for constructing a smooth mapping  $F(x, y)$  generally contains the following components:

- (1) An algorithm for forming a triangular grid for the convex hull of  $\{(x_i, y_i) : i = 1, \dots, n\}$ . The data points  $(x_i, y_i)$  are the vertices of the triangular cells of the grid.
- (2) A procedure for estimating the first (and possibly higher order) partial derivatives of  $F(x, y)$  at the data points  $(x_i, y_i)$ . There is currently no known best method for performing this task. At a point  $(x_i, y_i)$ , it is agreed that the derivative estimation should depend not only on  $z_i$ , but also on  $z_j$  for neighboring points  $(x_j, y_j)$ . However, the number of neighboring points that should be used in the derivative estimation is normally unclear.
- (3) For any point  $(x, y)$  in the grid, a routine for finding the triangular cell which contains the point. If extrapolation is to be permitted, then the region outside of the grid must be partitioned and a routine provided for locating any point which lies outside the grid.
- (4) A smooth interpolating algorithm for evaluating  $F(x, y)$  on each triangular cell of the grid. If extrapolation is to be permitted, then an algorithm must also be provided to compute  $F(x, y)$  on each cell of the partitioned region outside of the grid.

Generally, the derivative estimation appears to be the most ad hoc portion of most smooth interpolating procedures. This quite probably is unavoidable at the present time. However, it is unfortunate since the manner in which the derivative estimation is performed can significantly affect the results obtained from any interpolating procedure.

The subroutines BVIP and BVIP2 are available for obtaining a continuously differentiable surface  $z = F(x, y)$  for which  $z_i = F(x_i, y_i)$  for  $i = 1, \dots, n$ . Extrapolation is allowed, and the user is permitted to specify (via the argument  $n_c$ ) the number of neighboring points to be used for derivative estimation. BVIP is used if  $F(x, y)$  is to be evaluated on an arbitrary collection of output points, whereas BVIP2 is applicable only if  $F(x, y)$  is to be evaluated on a rectangular grid of output points. If BVIP is employed to evaluate  $F(x, y)$  on a rectangular grid, then BVIP will produce the same results as BVIP2 but it will be less efficient.

**CALL BVIP(MO,  $n_c$ ,  $n$ , X, Y, Z,  $m$ , XI, YI, ZI, IWK, WK, IERR)**

$X$  is an array containing  $x_1, \dots, x_n$ ,  $Y$  an array containing  $y_1, \dots, y_n$ , and  $Z$  an array containing  $z_1, \dots, z_n$ . The input argument  $n_c$  is the number of neighboring points to be used for derivative estimation. It is assumed that  $2 \leq n_c < n$  and  $n_c \leq 25$ . Currently no theory is available for indicating how  $n_c$  should be set. The only comment that can be made is that setting  $n_c$  to 3, 4, or 5 normally produces satisfactory results.

It is assumed that  $F(x, y)$  is to be evaluated at the points  $(\bar{x}_1, \bar{y}_1), \dots, (\bar{x}_m, \bar{y}_m)$ . XI is an array containing  $\bar{x}_1, \dots, \bar{x}_m$ , YI an array containing  $\bar{y}_1, \dots, \bar{y}_m$ , and ZI an array of dimension  $m$  or larger. When BVIP is called, if no input errors are detected then  $F(\bar{x}_i, \bar{y}_i)$  is computed and stored in ZI(i) for  $i = 1, \dots, m$ .

One may wish to recall BVIP a number of times to compute  $F(x, y)$  for different sets of points. If recalls are needed then a portion of the information that is generated on the first call to BVIP can frequently be reused. The reuse of information is controlled by the input argument MO. MO must have the value 1 on the first call to BVIP. For subsequent calls MO may be assigned the following values:

MO = 1 This setting must be employed when any of the data  $n_c, n, X, Y$  is modified. In this case, none of the previously generated information can be reused.

MO = 2 This setting may be used when  $n_c, n, X, Y$  are not modified.

MO = 3 This setting is permissible when only Z is modified.

If  $MO \neq 1$  then the contents of IWK and WK must not be altered.

IWK is an array of dimension  $kn + m$  or larger where  $k = \max\{31, 27 + n_c\}$ , and WK is an array of dimension  $8n$  or larger. IWK and WK are storage areas for the routine.

**Error Return.** IERR is an integer variable. If no errors are detected then IERR is set to 0. Otherwise, IERR is assigned one of the following values:

IERR = 1 MO is not 1, 2, or 3.

IERR = 2 Either  $2 \leq n_c < n$  or  $n_c \leq 25$  is violated.

IERR = 3  $n < 4$

IERR = 4  $m < 1$

IERR = 5 Either  $n_c$  or  $n$  has been modified. This cannot occur if  $MO \neq 1$ .

IERR = 6 The argument  $m$  has been modified. This cannot occur if  $MO = 3$ .

IERR = 7 Points  $(x_i, y_i)$  and  $(x_j, y_j)$  are equal or too close where  $IWK(1) = i$  and  $IWK(2) = j$ .

IERR = 8 The points  $(x_i, y_i, z_i)$  ( $i = 1, \dots, n$ ) are collinear or almost collinear.

When an error is detected, the routine immediately terminates.

#### Remarks.

- (1) The procedure is invariant under a rotation of the  $x$ - $y$  coordinate system.
- (2) The results are exact when  $F(x, y)$  represents a plane.
- (3) Derivative estimation at a data point depends on points closest to the data point. Thus the procedure is dependent on the scaling of the abscissae  $x_i$  and ordinates  $y_i$  of the data points.

**Programming.** BVIP employs the subroutines IDTANG, IDCLDP, IDLCTN, IDPDRV, IDPTIP and the function IDXCHG. The routines save and exchange information in labeled common blocks. The block names are IDLC and IDPI. The routines were written by Hiroshi Akima (Institute for Telecommunication Sciences, Boulder, Colorado). IDPTIP was modified by Albrecht Preusser (Fritz-Haber-Institut der Max-Planck-Gesellschaft, Berlin). The error handling was modified by A. H. Morris.

## References.

- (1) Akima, Hiroshi, "A Method of Bivariate Interpolation and Smooth Surface Fitting for Irregularly Distributed Data Points," *ACM Trans. Math Software* 4 (1978), pp. 148-159.
- (2) Preusser, A., "Remark on Algorithm 526," *ACM Trans. Math Software* 11 (1985), pp. 186-187.

**CALL BVIP2(MO,  $n_c$ ,  $n$ ,  $X$ ,  $Y$ ,  $Z$ ,  $\ell$ ,  $m$ , XI, YI, ZI, IWK, WK, IERR)**

$X$  is an array containing  $x_1, \dots, x_n$ ,  $Y$  an array containing  $y_1, \dots, y_n$ , and  $Z$  an array containing  $z_1, \dots, z_n$ . The input argument  $n_c$  is the number of neighboring points to be used for derivative estimation. It is assumed that  $2 \leq n_c < n$  and  $n_c \leq 25$ . Currently no theory is available for indicating how  $n_c$  should be set. The only comment that can be made is that setting  $n_c$  to 3, 4, or 5 normally produces satisfactory results.

It is assumed that  $F(x, y)$  is to be evaluated at  $(\bar{x}_i, \bar{y}_j)$  for  $i = 1, \dots, \ell$  and  $j = 1, \dots, m$ . XI is an array containing  $\bar{x}_1, \dots, \bar{x}_\ell$ , YI an array containing  $\bar{y}_1, \dots, \bar{y}_m$ , and ZI a 2-dimensional array of dimension  $\ell \times m$ . When BVIP2 is called, if no input errors are detected then  $F(\bar{x}_i, \bar{y}_j)$  is computed and stored in  $ZI(i, j)$  for  $i = 1, \dots, \ell$  and  $j = 1, \dots, m$ .

One may wish to recall BVIP2 a number of times for different grids  $(\bar{x}_i, \bar{y}_j)$ . If results are needed then a portion of the information that is generated on the first call to BVIP2 can frequently be reused. The reuse of information is controlled by the input argument MO. MO must have the value 1 on the first call to BVIP2. For subsequent calls MO may be assigned the following values:

MO = 1 This setting must be employed when any of the data  $n_c, n, X, Y$  is modified. In this case, none of the previously generated information can be reused.

MO = 2 This setting may be used when  $n_c, n, X, Y$  are not modified.

MO = 3 This setting is permissible when only  $Z$  is modified.

If MO  $\neq$  1 then the contents of IWK and WK must not be altered.

IWK is an array of dimension  $kn + \ell m$  or larger when  $k = \max\{31, 27 + n_c\}$ , and WK is an array of dimension  $5n$  or larger. IWK and WK are storage areas for the routine.

**Error Return.** IERR is an integer variable. If no errors are detected then IERR is set to 0. Otherwise, IERR is assigned one of the following values:

IERR = 1 MO is not 1, 2, or 3.

IERR = 2 Either  $2 \leq n_c < n$  or  $n_c \leq 25$  is violated.

IERR = 3  $n < 4$ .

IERR = 4 Either  $\ell < 1$  or  $m < 1$ .

IERR = 5 Either  $n_c$  or  $n$  has been modified. This cannot occur if MO  $\neq$  1.

IERR = 6 Either  $\ell$  or  $m$  has been modified. This cannot occur if MO = 3.

IERR = 7 Points  $(x_i, y_i)$  and  $(x_j, y_j)$  are equal or too close where  $IWK(1)=i$  and  $IWK(2) = j$ .

IERR = 8 The points  $(x_i, y_i, z_i)$  ( $i = 1, \dots, n$ ) are collinear or almost collinear.

When an error is detected, the routine immediately terminates.

**Remarks.**

- (1) The procedure is invariant under a rotation of the  $x$ - $y$  coordinate system.
- (2) The results are exact when  $F(x, y)$  represents a plane.
- (3) Derivative estimation at a data point depends on points closest to the data point. Thus the procedure is dependent on the scaling of the abscissae  $x_i$  and ordinates  $y_i$ .

**Programming.** BVIP2 employs the subroutines IDTANG, IDCLDP, IDGRID, IDPDRV, IDPTIP and the function IDXCHG. The routines save and exchange information in a labeled common block named IDPI. The routines were written by Hiroshi Akima (Institute for Telecommunication Sciences, Boulder, Colorado). IDPTIP was modified by Albrecht Preusser (Fritz-Haber-Institut der Max-Planck-Gesellschaft, Berlin). The error handling was modified by A. H. Morris.

**References.**

- (1) Akima, Hiroshi, "A Method of Bivariate Interpolation and Smooth Surface Fitting for Irregularly Distributed Data Points," *ACM Trans. Math Software* 4 (1978), pp. 148-159.
- (2) Preusser, A., "Remark on Algorithm 526," *ACM Trans. Math Software* 11 (1985), pp. 186-187.

## WEIGHTED LEAST SQUARES FITTING WITH POLYNOMIALS OF N VARIABLES

Let  $\{(x_1^{(i)}, \dots, x_n^{(i)}) : i = 1, \dots, \ell\}$  be a set of  $\ell$  distinct points,  $z_1, \dots, z_\ell$  be the corresponding function values to be approximated, and  $w_1, \dots, w_\ell$  be positive weights. Then for any nonnegative integer IDEG where  $\binom{n+\text{IDEG}}{n} \leq \ell$ ,<sup>1</sup> the subroutines MFIT and DMFIT are available for obtaining the coefficients of the unique polynomial  $F(x_1, \dots, x_n)$  of degree IDEG which minimizes  $\sum_{i=1}^{\ell} w_i [F(x_1^{(i)}, \dots, x_n^{(i)}) - z_i]^2$ . Also, the subroutines MEVAL and DMEVAL are available for computing this polynomial. MFIT and MEVAL yield single precision results, and DMFIT and DMEVAL yield double precision results.

```
CALL MFIT(n, IDEG, m, l, X, kx, Z, W, R, IER, IWK, WK, LIWK, LWK,
          MIWK, MWK)
CALL DMFIT(n, IDEG, m, l, X, kx, Z, W, R, IER, IWK, WK, LIWK, LWK,
           MIWK, MWK)
```

It is assumed that  $n \geq 1$  and  $\ell \geq 1$ .  $X$  is an  $\ell \times n$  matrix whose  $i^{\text{th}}$  row contains the point  $(x_1^{(i)}, \dots, x_n^{(i)})$ ; i.e.,  $X(i, j) = x_j^{(i)}$  for  $i = 1, \dots, \ell$  and  $j = 1, \dots, n$ . The argument  $kx$  is the number of rows in the dimension statement for  $X$  in the calling program.  $Z$  is an array containing  $z_1, \dots, z_\ell$  and  $W$  an array containing  $w_1, \dots, w_\ell$ .  $X$  and  $Z$  are modified by the routine.

**Remark.** For  $\text{IDEG} \geq 0$ ,  $\binom{n+\text{IDEG}}{n}$  polynomials  $1, x_1, \dots, x_n, x_1^2, x_1x_2, \dots$  are needed for a basis of the space of polynomials of degree  $\leq \text{IDEG}$ . The basis polynomials are ordered. For  $k \geq 1$ , the degree  $k-1$  basis polynomials precede the degree  $k$  polynomials. The degree  $k$  basis polynomials are  $x_{i_1} \cdots x_{i_k}$  where  $1 \leq i_1 \leq \dots \leq i_k \leq n$ . For any two such polynomials  $x_{i_1} \cdots x_{i_k}$  and  $x_{j_1} \cdots x_{j_k}$ , let  $r$  be the smallest integer such that  $i_r \neq j_r$ . Then  $x_{i_1} \cdots x_{i_k}$  precedes  $x_{j_1} \cdots x_{j_k}$  when  $i_r < j_r$ .

IDEG and  $m$  are variables. If  $\text{IDEG} \geq 0$  then the routine attempts to obtain the polynomial  $F(x_1, \dots, x_n)$  of degree IDEG which is the best least squares fit. Otherwise, if  $\text{IDEG} < 0$  then it is assumed that  $m \geq 1$  and that the first  $m$  basis polynomials are to be used to obtain the least squares fit. When the routine terminates, IDEG = the degree of the polynomial  $F(x_1, \dots, x_n)$  obtained and  $m$  = the number of basis polynomials that are actually used.

$R$  is an array of dimension  $\ell$  or larger.  $R(i) = z_i - F(x_1^{(i)}, \dots, x_n^{(i)})$  for  $i = 1, \dots, \ell$  when the routine terminates.

IWK is an array of dimension LIWK and WK an array of dimension LWK. When the routine terminates, IWK and WK contain the information needed for computing the polynomial  $F(x_1, \dots, x_n)$ . Sufficient storage for IWK and WK can be assured by setting LIWK and LWK as follows: If  $\text{IDEG} \geq 0$  then let  $N = \min\{\ell, \binom{n+\text{IDEG}}{n}\}$  and  $\delta = \text{IDEG}$ .

$$1 \binom{k}{0} = 1 \text{ and } \binom{k}{i} = \frac{k(k-1) \cdots (k-i+1)}{i!} \text{ for } i = 1, 2, \dots$$

Otherwise, if IDEG < 0 then let  $N = m$  and  $\delta$  be the smallest nonnegative integer such that  $\binom{n+\delta}{n} \geq m$ . Then set

$$\text{LIWK} \geq 4N + \delta n$$

$$\text{LWK} \geq 2N + n + 1 + \ell N + \frac{1}{2}(N-1)(N-2).$$

$N$  is the maximum number of basis polynomials that will be used, and  $\delta$  is the degree of the polynomial  $F(x_1, \dots, x_n)$  if  $N$  basis polynomials are used.

MIWK and MWK are variables. MIWK is set by the routine to the dimension needed for IWK, and MWK is set to the dimension needed for WK. MIWK and MWK depend only on  $n, \ell, \text{IDEG}$ , and  $m$ .

If MFIT is called then  $X, Z, W, R$ , and WK are single precision real arrays. Otherwise, if DMFIT is called then  $X, Z, W, R$ , and WK are double precision arrays.

IER is a variable that reports the status of the results. When the routine terminates, IER has one of the following values:

- IER = 0 The desired polynomial was obtained.
- IER = -1 Not all the basis polynomials could be used. IDEG is the degree of the polynomial obtained. This setting occurs when the problem is not solvable or is too ill-conditioned for the requested degree.
- IER = 1 Only  $\ell$  basis polynomials were used. A polynomial  $F(x_1, \dots, x_n)$  was obtained which solves the equations  $F(x_1^{(i)}, \dots, x_n^{(i)}) = z_i$  for  $i = 1, \dots, \ell$ .
- IER = 2 (Input error) IDEG < 0 and  $m \leq 0$ .
- IER = 3 (Input error)  $n < 1$  or  $\ell < 1$ .
- IER = 4 (Input error) LIWK or LWK is too small. Set  $\text{LIWK} \geq \text{MIWK}$  and  $\text{LWK} \geq \text{MWK}$ .

When an input error is detected, the routine immediately terminates.

**Remark.** When  $\text{IER} \leq 1$  then MEVAL or DMEVAL may be used to compute the polynomial obtained.

**Algorithm.** The revised Gram-Schmidt orthogonalization procedure is used.

**Programming.** MFIT employs the routines ALLOT, BASIZ, MTABLE, GNRTP, INCDG, SCALPM, SCALDN, and DMFIT employs the routines ALLOT, BASIZ, MTABLE, DGNRTP, DINCDG, DSCALP, DSCALD. MFIT and DMFIT are modifications by A. H. Morris of CONSTR, written by Richard H. Bartels (University of Waterloo) and John J. Jezioranski (Ontario Cancer Institute).

#### References.

- (1) Bartels, R. H. and Jezioranski, J. J., "Least Squares Fitting using Orthogonal Multinomials," *ACM Trans. Math Software* 11 (1985), pp. 201-217.
- (2) ———, "Algorithm 634, CONSTR and EVAL: Routines for Fitting Multinomials in a Least Squares Sense," *ACM Trans. Math Software* 11 (1985), pp. 218-228.

```

CALL MEVAL(n,KDEG, $\bar{m}$ , $\bar{l}$ ,XI,kxi,ZI,IND,IWK,WK,LIWK,LWK,T)
CALL DMEVAL(n,KDEG, $\bar{m}$ , $\bar{l}$ ,XI,kxi,ZI,IND,IWK,WK,LIWK,LWK,T)

```

MEVAL (DMEVAL) computes the polynomial obtained by MFIT (DMFIT), or a portion thereof. Let IDEG and  $m$  be the output values given by MFIT (DMFIT).

The argument  $\bar{m}$  is a variable. If  $KDEG < 0$  then it is assumed that  $1 \leq \bar{m} \leq m$  and that the polynomial using the first  $\bar{m}$  basis polynomials is to be computed. In this case, the polynomial computed is the best least squares fit for the basis polynomials involved.

If  $KDEG \geq 0$  then it is assumed that  $KDEG \leq IDEG$ . In this case, when the routine terminates,  $\bar{m} =$  the number of basis polynomials used. If  $\bar{m} \leq m$  (which will be the case when  $KDEG < IDEG$ ), then the polynomial computed is the polynomial of degree KDEG which is the best least squares fit.

**Usage.** If  $IER = \pm 1$  when MFIT(DMFIT) terminates, then the setting  $KDEG = IDEG$  normally causes an error to occur since  $\bar{m} > m$ . Hence, if it is desired that the full polynomial obtained by MFIT(DMFIT) be computed, no matter whether the value for  $IER$  is 0 or  $\pm 1$ , then KDEG should be assigned a negative value and  $\bar{m} = m$ .

It is assumed that the polynomial is to be computed at the points  $(\bar{x}_1^{(i)}, \dots, \bar{x}_n^{(i)})$  for  $i = 1, \dots, \bar{l}$ . XI is an  $\bar{l} \times n$  matrix whose  $i^{th}$  row contains the point  $(\bar{x}_1^{(i)}, \dots, \bar{x}_n^{(i)})$ . The argument *kxi* is the number of rows in the dimension statement for XI in the calling program. ZI is an array of dimension  $\bar{l}$  or larger. When the routine terminates,  $ZI(i)$  contains the value of the polynomial at the point  $(\bar{x}_1^{(i)}, \dots, \bar{x}_n^{(i)})$  for  $i = 1, \dots, \bar{l}$ .

IWK and WK are the arrays obtained from MFIT or DMFIT. LIWK is the dimension of IWK and LWK the dimension of WK. T is an array of dimension  $n$  or larger that is a work space for the routine.

If MEVAL is called then XI, ZI, WK and T are single precision arrays. Otherwise, if DMEVAL is called then XI, ZI, WK and T are double precision arrays.

IND is a variable that reports the status of the results. When the routine terminates, IND has one of the following values:

- IND = 0 The desired computation was performed.
- IND = -1 (Input error)  $\bar{m} < 1$  or  $\bar{m} > m$ .
- IND = -2 (Input error)  $n < 1$  or  $\bar{l} < 1$ .

**Programming.** MEVAL calls the subroutine MEVAL1 and DMEVAL calls the subroutine DMEVL1. MEVAL and DMEVAL are modifications by A. H. Morris of EVAL, written by Richard H. Bartels (University of Waterloo) and John J. Jezioranski (Ontario Cancer Institute).



## EVALUATION OF INTEGRALS OVER FINITE INTERVALS

QAGS, QSUBA, and DQAGS are available for computing integrals over finite intervals. The subroutine QAGS and function QSUBA yield single precision results, and the subroutine DQAGS yields double precision results. These procedures are adaptive. In such procedures, the selection of the points at which the integrand is evaluated depends on the nature of the integrand.

**CALL QAGS**( $F, a, b, \text{AERR}, \text{RERR}, z, \text{ERROR}, \text{NUM}, \text{IERR}, \ell, m, n, \text{IWK}, \text{WK}$ )

$F(x)$  is a user defined function whose arguments and values are assumed to be single precision real numbers. The purpose of QAGS is to compute the integral  $\int_a^b F(x) dx$ .  $F$  need not be defined at  $a$  and  $b$ , and it is not required that  $a \leq b$ .  $F$  must be declared in the calling program to be of type EXTERNAL.

AERR and RERR are the absolute and relative error tolerances to be used, and  $z$  is a variable. When QAGS is called,  $z$  is assigned the value obtained for  $\int_a^b F(x) dx$ . The routine attempts to obtain a value  $z$  which satisfies  $|\int_a^b F(x) dx - z| \leq \max\{\text{AERR}, \text{RERR} \cdot |\int_a^b F(x) dx|\}$ . It is assumed that  $\text{AERR} \geq 0$  and  $\text{RERR} \geq 0$ . If one wants accuracy to  $k$  significant digits then set  $\text{RERR} = 10^{-k}$ .

ERROR and NUM are variables. When QAGS terminates, ERROR is the estimated absolute error of the result and NUM is the number of points at which  $F$  was evaluated.

IWK is an array of dimension  $\ell$  and WK is an array of dimension  $m$ . IWK and WK are work spaces for the routine. The input argument  $\ell$  is the maximum number of subintervals in which the interval of integration may be partitioned. It is assumed that  $\ell \geq 1$  and  $m \geq 4\ell$ . The argument  $n$  is a variable. When QAGS terminates,  $n$  = the number of subintervals that appeared in the partition. Normally  $n < 100$ .

IERR is a variable that reports the status of the results. When the routine terminates, IERR has one of the following values:

- IERR = 0 The routine is satisfied that the integral has been computed to the desired accuracy.
- IERR = 1 The interval of integration was partitioned into  $\ell$  subintervals. More subintervals are needed to compute the integral to the desired accuracy.
- IERR = 2 The integral has been computed, but because of roundoff error QAGS is not certain of the accuracy of the result. The error may be greater than that reported by ERROR.
- IERR = 3 Extremely bad integrand behavior occurs in the interval of integration. The routine is not certain of the accuracy obtained.
- IERR = 4 The algorithm does not converge. It is assumed that the requested accuracy cannot be achieved and that the result is the best which can be obtained.

IERR = 5 The integral may be divergent or it may converge extremely slowly.  
In this case, the value for  $z$  may be meaningless.

IERR = 6 (Input Error) Either  $\ell < 1$ ,  $m < 4\ell$ , AERR < 0, or RERR < 0. In  
this case, the variables  $z$ , ERROR, NUM, and  $n$  are set to 0.

**Note.**  $F$  may have singularities at  $a$  and  $b$ . However, it is recommended that no singularities appear in the interior of the interval of integration.

**Algorithm.** The 21 point Kronrod rule and  $\epsilon$ -algorithm of P. Wynn are used.

**Programming.** QAGS employs the subroutines QAGSE, QK21F, QPSRT, and QELG. These routines were developed by Robert Piessens and Elise de Doncker-Kapenga (Katholieke Universiteit Leuven, Heverlee, Belgium). The function SPMPAR is also used.

**Reference.** Piessens, R., de Doncker-Kapenga, E., Uberhuber, C. W., and Kahaner, D. K., *QUADPACK: A Subroutine Package for Automatic Integration*, Springer-Verlag, 1983.

**QSUBA**( $F, a, b, \text{RERR}, \text{MCOUNT}, \text{ERROR}, \text{IND}$ )

$F(x)$  is a user defined function whose arguments and values are assumed to be single precision real numbers. The purpose of QSUBA is to compute the integral  $\int_a^b F(x) dx$ .  $F$  need not be defined at the points  $a$  and  $b$ . However, it is assumed that  $a \leq b$ .  $F$  must be declared in the calling program to be of type EXTERNAL.

RERR is the relative error tolerance to be satisfied. It is assumed that  $\text{RERR} > 0$ . If one wants accuracy to  $k$  significant digits then set  $\text{RERR} = 10^{-k}$ .

The input argument MCOUNT is the maximum number of points at which  $F$  may be evaluated. It is recommended that  $\text{MCOUNT} \geq 1000$ .

ERROR is a variable that is set by QSUBA. If the value of QSUBA is not 0 then ERROR is a rough estimate of the relative error of the computed result. Otherwise, if the value of QSUBA is 0 then ERROR is rough estimated of the absolute error.

IND is a variable that reports the status of the results. When QSUBA terminates, IND has one of the following values:

IND = 0 The function QSUBA is satisfied that the integral has been computed to the desired accuracy.

IND = 1 The integral has been computed, but QSUBA is not certain of the accuracy of the result.

IND = 2  $F(x)$  was evaluated at MCOUNT points. More evaluations are needed to complete the computation of the integral.

IND = 3 The function QSUBA cannot compute the integral to the desired accuracy.

If  $\text{IND} = 0$  or  $1$  then the function QSUBA is assigned the value obtained for the integral. If  $\text{IND} = 2$  then QSUBA has for its value the most recent acceptable partial estimate made of the integral. Otherwise, if  $\text{IND} = 3$ , then QSUBA has for its value the best estimate of the value of the integral that it can make.

**Note.** QSUBA assumes that  $F$  and its derivatives have no singularities in the closed interval  $[a, b]$ . If this is not the case then QAGS should be used. QSUBA is recommended for computing integrals such as  $\int_0^{10} \sin \frac{\pi}{2} x^2 dx$  whose integrands are finitely oscillatory.

**Algorithm.** Gaussian quadrature is employed.

**Programming.** QSUBA calls the subroutine QUAD. QSUBA and QUAD were written by T. N. L. Patterson (Queen's University, Belfast, Northern Ireland), and QSUBA was modified by A. H. Morris. The function SPMPAR is used.

**Reference.** Patterson, T. N. L., "Algorithm for Automatic Numerical Integration Over a Finite Interval," *Comm. ACM* 16 (1973), pp.694-699.

**CALL DQAGS**( $F, a, b, AERR, RERR, z, ERROR, NUM, IERR, \ell, m, n, IWK, WK$ )

$F(x)$  is a user defined function whose arguments and values are assumed to be double precision real numbers. The purpose of DQAGS is to compute the integral  $\int_a^b F(x) dx$ . The arguments  $a$  and  $b$  are double precision real numbers.  $F$  need not be defined at  $a$  and  $b$ , and it is not required that  $a \leq b$ .  $F$  must be declared in the calling program to be of types DOUBLE PRECISION and EXTERNAL.

AERR and RERR are double precision real numbers and  $z$  is a double precision variable. AERR and RERR are the absolute and relative error tolerances to be used. When DQAGS is called,  $z$  is assigned the value obtained for  $\int_a^b F(x) dx$ . The routine attempts to obtain a value  $z$  which satisfies  $|\int F(x) dx - z| \leq \max\{AERR, RERR \cdot |\int F(x) dx|\}$ . It is assumed that AERR  $\geq 0$  and RERR  $\geq 0$ .

ERROR is a double precision variable and NUM an integer variable. When DQAGS terminates, ERROR is the estimated absolute error of the result and NUM is the number of points at which  $F$  was evaluated.

IWK is an integer array of dimension  $\ell$  and WK a double precision array of dimension  $m$ . IWK and WK are work spaces for the routine. The argument  $\ell$  is the maximum number of subintervals in which the interval of integration may be partitioned. It is assumed that  $\ell \geq 1$  and  $m \geq 4\ell$ . The argument  $n$  is a variable. When DQAGS terminates,  $n$  = the number of subintervals that appeared in the partition. Normally  $n < 100$ .

IERR is a variable that reports the status of the results. When the routine terminates, IERR has one of the following values:

- IERR = 0 The routine is satisfied that the integral has been computed to the desired accuracy.
- IERR = 1 The interval of integration was partitioned into  $\ell$  subintervals. More subintervals are needed to compute the integral to the desired accuracy.
- IERR = 2 The integral has been computed, but because of roundoff error DQAGS is not certain of the accuracy of the result. The error may be greater than that reported by ERROR.

- IERR = 3 Extremely bad integrand behavior occurs in the interval of integration. The routine is not certain of the accuracy obtained.
- IERR = 4 The algorithm does not converge. It is assumed that the requested accuracy cannot be achieved and that the result is the best which can be obtained.
- IERR = 5 The integral may be divergent or it may converge extremely slowly. In this case, the value for  $z$  may be meaningless.
- IERR = 6 (Input Error) Either  $\ell < 1$ ,  $m < 4\ell$ , AERR < 0, or RERR < 0. In this case, the variables  $z$ , ERROR, NUM, and  $n$  are set to 0.

**Remarks.**  $F$  may have singularities at  $a$  and  $b$ . However, it is recommended that no singularities appear in the interior of the interval of integration. DQAGS is a double precision version of the routine QAGS.

**Algorithm.** The 21 point Kronrod rule and  $\epsilon$ -algorithm of P. Wynn are used.

**Programming.** DQAGS employs the subroutines DQAGSE, DQK21, DQPSRT, and DQELG. These subroutines are double precision versions of the subroutines QAGSE, QK21F, QPSRT, and QELG, developed by Robert Piessens and Elise de Doncker-Kapenga (Katholieke Universiteit Leuven, Heverlee, Belgium). The function DPMPAR is also used.

**Reference.** Piessens, R., de Doncker-Kapenga, E., Uberhuber, C. W., and Kahaner, D. K., *QUADPACK: A Subroutine Package for Automatic Integration*, Springer-Verlag, 1983.

## EVALUATION OF INTEGRALS OVER INFINITE INTERVALS

The subroutines QAGI and DQAGI are available for computing integrals over infinite intervals. QAGI yields single precision results and DQAGI yields double precision results. QAGI and DQAGI are adaptive routines.

**CALL QAGI**( $F,a,MO,AERR,RERR,z,ERROR,NUM,IERR,\ell,m,n,IWK,WK$ )

$F(x)$  is a user defined function whose arguments and values are assumed to be real numbers. The argument  $a$  is a real number,  $z$  is a variable, and  $MO$  may be 1, -1, or 2. When QAGI is called,  $z$  is assigned the value  $\int_a^\infty F(x) dx$  if  $MO = 1$  and the value  $\int_{-\infty}^a F(x) dx$  if  $MO = -1$ . Otherwise, if  $MO = 2$  then  $z$  is assigned the value  $\int_{-\infty}^\infty F(x) dx$ . If  $MO = \pm 1$  then  $F$  need not be defined at  $a$ . Otherwise, if  $MO = 2$  then  $a$  is not used.  $F$  must be declared in the calling program to be of type EXTERNAL.

AERR and RERR are the absolute and relative error tolerances to be used. The subroutine attempts to obtain a value  $z$  which satisfies  $|\int F(x) dx - z| \leq \max\{AERR, RERR \cdot |\int F(x) dx|\}$ . It is assumed that  $AERR \geq 0$  and  $RERR \geq 0$ . If one wants accuracy to  $k$  significant digits then set  $RERR = 10^{-k}$ .

ERROR and NUM are variables. When QAGI terminates, ERROR is the estimated absolute error of the result and NUM is the number of points at which  $F$  was evaluated.

IWK is an array of dimension  $\ell$  and WK is an array of dimension  $m$ . IWK and WK are work spaces for the routine. The input argument  $\ell$  is the maximum number of subintervals in which the interval of integration may be partitioned. It is assumed that  $\ell \geq 1$  and  $m \geq 4\ell$ . The argument  $n$  is a variable. When QAGI terminates,  $n$  = the number of subintervals that appeared in the partition. Normally  $n < 100$ .

IERR is a variable that reports the status of the results. When the routine terminates, IERR has one of the following values:

- IERR = 0 The routine is satisfied that the integral has been computed to the desired accuracy.
- IERR = 1 The interval of integration was partitioned into  $\ell$  subintervals. More subintervals are needed to compute the integral to the desired accuracy.
- IERR = 2 The integral has been computed, but because of roundoff error QAGI is not certain of the accuracy of the result. The error may be greater than that reported by ERROR.
- IERR = 3 Extremely bad integrand behavior occurs in the interval of integration. The routine is not certain of the accuracy obtained.
- IERR = 4 The algorithm does not converge. It is assumed that the requested accuracy cannot be achieved and that the result is the best which can be obtained.
- IERR = 5 The integral may be divergent or it may converge extremely slowly. In this case, the value for  $z$  may be meaningless.

IERR = 6 (Input Error) Either  $\ell < 1$ ,  $m < 4\ell$ , AERR < 0, or RERR < 0. In this case, the variables  $z$ , ERROR, NUM, and  $n$  are set to 0.

**Note.**  $F$  may have a singularity at  $a$  when  $MO = \pm 1$ . However, it is recommended that no singularities appear in the interior of the interval of integration.

**Algorithm.** The integrals are transformed as follows:

$$\begin{aligned}\int_a^\infty F(x) dx &= \int_0^1 F(a - 1 + 1/t) \frac{dt}{t^2} \\ \int_{-\infty}^a F(x) dx &= \int_0^1 F(a + 1 - 1/t) \frac{dt}{t^2} \\ \int_{-\infty}^\infty F(x) dx &= \int_0^1 [F(-1 + 1/t) + F(1 - 1/t)] \frac{dt}{t^2}\end{aligned}$$

The transformed integrals are computed by the 15 point Kronrod rule and the  $\epsilon$ -algorithm of P. Wynn.

**Programming.** QAGI employs the subroutines QAGIE, QK15I, QPSRT, and QELG. These routines were developed by Robert Piessens and Elise de Doncker-Kapenga (Katholieke Universiteit Leuven, Heverlee, Belgium). The function SPMPAR is also used.

**Reference.** Piessens, R., de Doncker-Kapenga, E., Uberhuber, C. W., and Kahaner, D. K., *QUADPACK: A Subroutine Package for Automatic Integration*, Springer-Verlag, 1983.

**CALL DQAGI**( $F, a, MO, AERR, RERR, z, ERROR, NUM, IERR, \ell, m, n, IWK, WK$ )

$F(x)$  is a user defined function whose arguments and values are assumed to be double precision real numbers. The argument  $a$  is a double precision real number,  $z$  is a double precision variable, and  $MO$  may be 1, -1, 2. When DQAGI is called,  $z$  is assigned the value  $\int_a^\infty F(x) dx$  if  $MO = 1$  and the value  $\int_{-\infty}^a F(x) dx$  if  $MO = -1$ . Otherwise, if  $MO = 2$  then  $z$  is assigned the value  $\int_{-\infty}^\infty F(x) dx$ . If  $MO = \pm 1$  then  $F$  need not be defined at  $a$ .  $F$  must be declared in the calling program to be of types DOUBLE PRECISION and EXTERNAL.

AERR and RERR are the absolute and relative error tolerances to be used. The subroutine attempts to obtain a value  $z$  which satisfies  $|\int F(x) dx - z| \leq \max\{AERR, RERR \cdot |\int F(x) dx|\}$ . It is assumed that AERR and RERR are nonnegative double precision numbers.

ERROR is a double precision variable and NUM an integer variable. When DQAGI terminates, ERROR is the estimated absolute error of the result and NUM is the number of points at which  $F$  was evaluated.

IWK is an integer array of dimension  $\ell$  and WK a double precision array of dimension  $m$ . IWK and WK are work spaces for the routine. The argument  $\ell$  is the maximum number of subintervals in which the interval of integration may be partitioned. It is assumed that  $\ell \geq 1$  and  $m \geq 4\ell$ . The argument  $n$  is a variable. When DQAGI terminates,  $n$  = the number of subintervals that appeared in the partition. Normally  $n < 100$ .

IERR is a variable that reports the status of the results. When the routine terminates, IERR has one of the following values:

- IERR = 0 The routine is satisfied that the integral has been computed to the desired accuracy.
- IERR = 1 The interval of integration was partitioned into  $\ell$  subintervals. More subintervals are needed to compute the integral to the desired accuracy.
- IERR = 2 The integral has been computed, but because of roundoff error DQAGI is not certain of the accuracy of the result. The error may be greater than that reported by ERROR.
- IERR = 3 Extremely bad integrand behavior occurs in the interval of integration. The routine is not certain of the accuracy obtained.
- IERR = 4 The algorithm does not converge. It is assumed that the requested accuracy cannot be achieved and that the result is the best which can be obtained.
- IERR = 5 The integral may be divergent or it may converge extremely slowly. In this case, the value for  $z$  may be meaningless.
- IERR = 6 (Input Error) Either  $\ell < 1$ ,  $m < 4\ell$ , AERR  $< 0$ , or RERR  $< 0$ . In this case, the variables  $z$ , ERROR, NUM, and  $n$  are set to 0.

**Remarks.**  $F$  may have a singularity at  $a$  when  $MO = \pm 1$ . However, it is recommended that no singularities appear in the interior of the interval of integration. DQAGI is a double precision version of the routine QAGI.

**Programming.** DQAGI employs the routines DQAGIE, DQK15I, DQPSRT, and DQELG. These subroutines are double precision versions of the subroutines QAGIE, QK15I, QPSRT, and QELG, developed by Robert Piessens and Elise de Doncker-Kapenga (Katholieke Universiteit Leuven, Heverlee, Belgium). The function DPMPAR is also used.

**Reference.** Piessens, R., de Doncker-Kapenga, E., Uberhuber, C. W., and Kahaner, D. K., *QUADPACK: A Subroutine Package for Automatic Integration*, Springer-Verlag, 1983.

## EVALUATION OF DOUBLE INTEGRALS OVER TRIANGLES

Let  $f(x, y)$  be a real-valued function defined on a triangle  $T$ . Then the subroutine CUBTRI is available for computing the integral  $\iint_T f(x, y) dx dy$ . CUBTRI is an adaptive routine.

**CALL CUBTRI( $F, T, \epsilon, \text{MAX}, A, \text{ERR}, n, W, \ell, \text{IDATA}, \text{RDATA}, \text{IERR}$ )**

$T$  is a 2-dimensional real array of dimension  $2 \times 3$  where  $T(1, j)$  and  $T(2, j)$  are the  $x$  and  $y$  coordinates of the  $j^{\text{th}}$  vertex of the given triangle ( $j = 1, 2, 3$ ).

IDATA and RDATA are arrays provided by the user containing any integer or real data needed for computing the integrand  $f(x, y)$ . The arrays may be of any size.  $F$  is a user defined real-valued function having the arguments  $x, y, \text{IDATA}, \text{RDATA}$ . It is assumed that  $F(x, y, \text{IDATA}, \text{RDATA}) = f(x, y)$  for any point  $(x, y)$  in the triangle of integration  $T$ .  $F$  must be declared in the calling program to be of type EXTERNAL.

The input argument  $\epsilon$  is the error tolerance to be satisfied, and  $A$  is a variable. When CUBTRI is called,  $A$  is assigned the value obtained for  $\iint_T f(x, y) dx dy$ . The routine attempts to obtain a value  $A$  which satisfies  $|\iint f(x, y) dx dy - A| < \max\{\epsilon, \epsilon|A|\}$ . Thus if  $|A| \leq 1$  then  $\epsilon$  is an absolute tolerance, whereas if  $|A| > 1$  then  $\epsilon$  is a relative tolerance. If one wants  $k$  digit accuracy then set  $\epsilon = 10^{-k}$ . ERR is a variable. When CUBTRI terminates, ERR is the estimated error  $|\iint f(x, y) dx dy - A|$  of the result.

The input argument MAX is the maximum number of points  $(x, y)$  at which  $F$  may be evaluated, and  $n$  is a variable. On an initial call to CUBTRI, the user must set  $n = 0$ . When the routine terminates,  $n$  will have for its value the number of points at which  $F$  was evaluated. (For subsequent calls concerning the same integral, see below.)

$W$  is an array of dimension  $\ell$  for internal use by the routine. The input argument  $\ell$  specifies the maximum number of subtriangles in which the triangle of integration  $T$  may be partitioned. The subtriangles are stored in  $W$ , each subtriangle requiring 6 storage locations. Thus  $\ell/6$  is an estimate of the maximum number of subtriangles that might have to be stored ( $\ell \leq \max\{1, 3m + 1\}$  where  $m = (\text{MAX}/19 - 1)/4$ ).

IERR is an integer variable that reports the status of the results. When the routine terminates, IERR has one of the following values:

- IERR = 0 The integral was computed to the desired accuracy.
- IERR = 1 MAX is too small.  $F$  must be evaluated at more points.
- IERR = 2 The storage space  $W$  is full. Its dimension  $\ell$  must be increased.
- IERR = 3 Further subdivision of the subtriangles impossible. This normally occurs when  $f(x, y)$  has a singularity in the region. The situation can frequently be eliminated by placing the singularity at a vertex of the triangle of integration  $T$ .
- IERR = 4 No further improvement in accuracy is possible because of roundoff error in the computation of  $F$  or the irregular behavior of  $F$ .



IERR = 5 No further improvement in accuracy is possible because subdivision does not change the estimated integral value  $A$ . Machine accuracy has probably been reached.

After an initial call to CUBTRI, the routine may be recalled to continue the computation of  $\iint_T f(x, y) dx dy$ . When the routine is recalled, the value of  $n$  obtained on the previous call to CUBTRI is used for the next call. This value for  $n$  tells the routine where computation should be resumed (using the information previously stored in  $W$ ). At least one of the values  $\epsilon$ , MAX, or  $\ell$  must be modified before CUBTRI is recalled.  $F$ ,  $T$ ,  $n$ ,  $W$ , IDATA, and RDATA may not be changed when the routine is recalled.

**Remark.**  $F$  may have a singularity at one of the vertices of  $T$  (such as in the case when we are computing  $\int_0^1 \int_0^x (x^2 + 3y^2)^{-1/2} dy dx$ ). However, it is recommended that no singularities appear in the interior of the triangle of integration.

**Algorithm.** The 7-point degree 5 rule of Radon and a new 19-point degree 8 rule are used.

**Programming.** CUBTRI calls the function RNDERR and subroutine CUBRUL. Information is saved in labeled common blocks. The block names are CUBSTA and CUBATB. The routines were written by D. P. Laurie (National Research Institute for the Mathematical Sciences, Pretoria, South Africa).

**Reference.** Laurie, D. P., "Algorithm 584, CUBTRI: Automatic Cubature over a Triangle," *ACM Trans. Math Software* 8 (1982), pp. 210-218.

## SOLUTION OF FREDHOLM INTEGRAL EQUATIONS OF THE SECOND KIND

If  $k(s, t)$  and  $f(s)$  are continuous real-valued functions for  $a \leq s \leq b$  and  $a \leq t \leq b$ , then the equation to be solved is

$$x(s) - \int_a^b k(s, t)x(t) dt = f(s)$$

for  $a \leq s \leq b$ . Let  $K$  be the operator defined by  $(Kx)(s) = \int_a^b k(s, t)x(t) dt$  for any real-valued function  $x$  continuous on  $[a, b]$ . Then  $(Kx)(s)$  is continuous for  $a \leq s \leq b$ , and  $k$  is called the *kernel* of  $K$ . Also the above integral equation can be written in the form  $(I - K)x = f$  where  $I$  is the identity operator. This equation has a unique solution if and only if  $I - K$  is 1 - 1, in which case  $x = (I - K)^{-1}f$ . The subroutine IESLV is available for computing this solution.

**Remark.** If  $C[a, b]$  is the normed space of real-valued functions  $x$  continuous on  $[a, b]$  and having the norm  $\|x\| = \max\{|x(t)| : a \leq t \leq b\}$ , then  $K$  is a compact mapping  $C[a, b] \rightarrow C[a, b]$  having the norm  $\|K\| = \max_{a \leq s \leq b} \int_a^b |k(s, t)| dt$ .

**CALL IESLV(*k, f, a, b, EPS, IFLAG, S, X, l, N, M, NF, MF, NORMK, WK, IERR*)**

It is assumed that  $a < b$ , and that  $k(s, t)$  and  $f(s)$  are user defined real-valued functions for  $a \leq s, t \leq b$ . It is recommended that  $k$  and  $f$  be several times continuously differentiable. The functions  $k$  and  $f$  must be declared in the calling program to be of type EXTERNAL.

EPS is a variable and IFLAG an input argument whose values are 0 and 1. On input EPS is the error tolerance that the solution must satisfy. If IFLAG = 0 then EPS is an absolute tolerance. Otherwise, if IFLAG = 1 then EPS is a relative tolerance. If IESLV solves the equation, then on output EPS is the estimated error of the result.

Before the remaining arguments  $s, x, l, \dots$  can be described, it is necessary to give a brief outline of the algorithm used. When IESLV is called, the integral equation is approximated by

$$(*) \quad x_n(s) - \sum_{j=1}^n w_{jn} k(s, t_{jn}) x_n(t_{jn}) = f(s)$$

for  $a \leq s \leq b$ . Here  $w_{jn}$  and  $t_{jn}$  are the weights and nodes of Gauss-Legendre quadrature. This equation is treated as an interpolation for  $x(s)$  in terms of the values  $x_n(t_{jn})$ . These values are obtained by solving the equations

$$(**) \quad x_n(t_{in}) - \sum_{j=1}^n w_{jn} k(t_{in}, t_{jn}) x_n(t_{jn}) = f(t_{in})$$

for  $i = 1, \dots, n$ . This system of equations can be solved directly or iteratively. The following algorithm is used:

- (1) Set  $n = 2$  and go to (2).
- (2) The  $n$  equations are solved directly. Then set  $m = 2n$  and solve the  $m$  equations (\*\*) iteratively. If the rate of convergence is sufficiently rapid or  $n$  cannot be increased, then go to (3). Otherwise, set  $n = m$ , and go to (2).
- (3) Here  $n$  remains fixed. Repeatedly double the value of  $m$  and solve the  $m$  equations (\*\*) iteratively until convergence occurs,  $m$  cannot be increased, or the iterations diverge.

When the algorithm terminates, values  $x_m(t_{im})$  will have been computed for the nodes  $t_{im} (i = 1, \dots, m)$ . Then from (\*),  $x(s) \approx x_m(s)$  can be interpolated for  $a \leq s \leq b$ .

$N$  and  $M$  are input arguments, and  $WK$  is an array that is a work space for the routine.  $N$  and  $M$  are upper limits for  $n$  and  $m$  in the algorithm, and  $WK$  is of dimension  $5N^2 + 9(N + M)$  or larger. It is assumed that  $M \geq N \geq 2$ . Since  $n$  and  $m$  are always powers of 2,  $N$  and  $M$  need only be set to powers of 2. However, this is not required.

$S$  and  $X$  are arrays, and  $\ell$  is a variable. On input it is assumed that  $\ell \geq 0$ . If  $\ell > 0$  then  $S$  is assumed to contain  $\ell$  points  $s_1, \dots, s_\ell$  at which the solution  $x(s)$  is to be evaluated. Also  $X$  is assumed to be an array of dimension  $\ell$  or larger. When IESLV terminates,  $X$  contains the values obtained for  $x(s_1), \dots, x(s_\ell)$ . (This is true irregardless of whether or not the desired accuracy has been achieved.) Otherwise, if  $\ell = 0$  then  $S$  and  $X$  are assumed to be arrays of dimension  $M$  or larger. When IESLV terminates  $\ell =$  the final value obtained for  $m$ ,  $S$  contains the Gaussian nodes  $t_{i\ell} (i = 1, \dots, \ell)$ , and  $X$  contains the values obtained for  $x(t_{i\ell})$ .

$NF$  and  $MF$  are variables. When the routine terminates,  $NF$  is the final value for  $n$  and  $MF$  the final value for  $m$ .

$NORMK$  is a real variable. If  $\ell > 0$  on input, then when IESLV terminates,  $NORMK$  is an approximation for  $\|K\|$ . Otherwise, if  $\ell = 0$  then  $NORMK = 0$ .

$IERR$  is a variable that reports the status of the results. When the routine terminates,  $IERR$  has one of the following values:

- $IERR = 0$  The solution was obtained to the desired accuracy.  $EPS$  is the estimated error of the result.
- $IERR = 1$  The solution was not obtained to the desired accuracy.  $EPS$  is the estimated error of the result.
- $IERR = 2$  The solution was not obtained to the desired accuracy. It is not clear what accuracy (if any) has been achieved.  $EPS$  has been set to 0.
- $IERR = 3$  The input value for  $EPS$  was too small. This may be due to ill-conditioning of the integral equation. The value of  $EPS$  was reset to a more realistic tolerance, which the solution satisfied.
- $IERR = 4$  The solution  $x(s)$  was obtained at the Gaussian nodes to the desired precision. However, the interpolation process may not preserve this accuracy for the evaluation of  $x(s)$  for other points  $s$ .  $EPS$  is the estimated error of the solution at the Gaussian nodes.

IERR = 5 The solution  $x(s)$  was not obtained to the desired accuracy at the Gaussian nodes. EPS is the estimated error at these nodes. The interpolation process may not preserve this accuracy for the evaluation of  $x(s)$  for other points  $s$ .

IERR = 6 The input value for EPS was too small. This may be due to ill-conditioning of the integral equation. The value of EPS was reset to a more realistic tolerance, which the solution  $x(s)$  satisfied at the Gaussian nodes. The interpolation process may not preserve this accuracy for the evaluation of  $x(s)$  for other points  $s$ .

Difficulties can arise, causing  $IERR \geq 1$ , when the integral equation is ill-conditioned or the kernel  $k(s, t)$  is not appropriate for standard Gaussian quadrature. Ill-conditioning can occur when the operator  $I - K$  is near singular or the norm  $\|K\|$  is exceedingly large. Inappropriate kernels  $k(s, t)$  include those which are highly oscillatory or not continuously differentiable for  $s$  and  $t$  in the open interval  $(a, b)$ .

**Programming.** IESLV employs the subroutines IEGS, NSTERP, WANDT, LEAVE, ITERT, LNSYS and functions RMIN, RNRM, CONEW. The routines save and exchange information in labeled common blocks. The block names are XXINFO and XXLIN. The routines were written by Kendall E. Atkinson (University of Iowa), and modified by A. H. Morris. The function SPMPAR is also used.

**Reference.** Atkinson, K. E., "An Automatic Program for Linear Fredholm Integral Equations of the Second Kind," *ACM Trans. Math Software* 2 (1976), pp. 154-171.

## THE INITIAL VALUE SOLVERS – INTRODUCTORY COMMENTS

Let  $y'(t) = f(t, y(t))$  denote a system of  $n$  ordinary first order differential equations where  $f(t, y) = (f_1(t, y), \dots, f_n(t, y))$  and  $y(t) = (y_1(t), \dots, y_n(t))$ . Assume that  $y(a)$  is known. Then for  $b \neq a$  the subroutines ODE, RKF45, GERK, SFODE, and SFODE1 are available for computing  $y(b)$ . These routines are adaptive variable step differential equations solvers. The remaining subroutines (RK and RK8) are fixed order, one step procedures. Given a value  $y(t)$  and a step size  $h$ , the one step procedures compute a value for  $y(t + h)$ . The problem of selecting an appropriate step size is left to the user. Given  $y(a)$  and  $b$ , the one step routines must be repeatedly called to step along the interval from  $a$  to  $b$ . The situation, however, is considerably different with the adaptive routines. ODE, SFODE, and SFODE1 are variable order, variable step procedures, and RKF45 and GERK are fixed order, variable step procedures. Given  $y(a)$ ,  $b$ , and the error tolerances that are to be maintained, these solvers continually adjust their orders and step sizes as they (automatically) step along the interval from  $a$  to  $b$ .

The adaptive routines differ in their capabilities. ODE, RKF45, and GERK are recommended for nonstiff equations, and SFODE and SFODE1 for stiff equations. If one does not know whether the equations are stiff, then ODE should be tried. ODE maintains greater accuracy than the other routines, and it will notify the user if the equations appear to be stiff. ODE, RKF45, and GERK should be able to handle mildly stiff problems satisfactorily, but they are decidedly not appropriate for extremely stiff problems. SFODE and SFODE1 are the only routines in the mathematics library that are capable of solving extremely stiff equations.

If the equations to be solved are nonstiff, then the choice between ODE and RKF45 depends on the amount of accuracy needed and the cost of the derivative evaluations. If the accuracy requirements are high then ODE is recommended. However, if the accuracy requirements are low and the derivative evaluations are inexpensive, then RKF45 may be the most efficient routine for the problem. RKF45 frequently requires more derivative evaluations than ODE, but its overhead is considerably less than that for ODE.

When the user specifies the error tolerances to be satisfied, normally he is only interested in the global error (the accuracy of  $y(b)$ ). However, the adaptive routines employ the tolerances for controlling local error (the error generated at each internal step in the interval). No attempt is made to control the progressive erosion of accuracy that can occur when the steps accumulate. GERK is the only routine that estimates the global error. This routine employs the same Runge-Kutta-Fehlberg formulae used by RKF45. GERK is 2-3 times slower than RKF45, but it is more accurate.

**Output Considerations.** Generally, when the user has a system of equations  $y'(t) = f(t, y(t))$  to be solved (where  $y(a_0)$  is known), he wants its solution at a sequence of points  $a_1, \dots, a_N$ . If an adaptive routine is being used, then the routine will be repeatedly called to step along the axis from each point  $a_i$  to the next. If ODE, SFODE, or SFODE1 is being employed, then the number and closeness of the output points  $a_1, \dots, a_N$  should be of no concern. These routines partially ignore  $a_{i+1}$  in the selection of the step size when

going from  $a_i$  to  $a_{i+1}$ . Instead, they step along the axis using the largest steps that are appropriate (efficiency and accuracy are the prime concerns). Normally  $a_{i+1}$  is passed in the process. If  $a_{i+1}$  is passed then a quick interpolation yields the desired result at  $a_{i+1}$ . Thus the process of solving the equations for  $a_{i+1}$  when  $y(a_i)$  is known may require that no steps be taken ( $a_{i+1}$  may have been bypassed on a previous call to ODE, SFODE, or SFODE1), or it may require that one or more steps be taken.

The situation is considerably different if RKF45 or GERK is used. These routines select their step size so as not to bypass  $a_{i+1}$  when going from  $a_i$  to  $a_{i+1}$ . Thus the output points  $a_1, \dots, a_N$  may be so close to one another as to force inordinately small step sizes (when such step sizes would otherwise not be needed). If this occurs then the efficiency of RKF45 and GERK may deteriorate dramatically. The routines will notify the user of the situation, and the user will be left with the following options:

- (1) Switch to an adaptive routine such as ODE which performs interpolation.
- (2) Use a nonadaptive one step routine such as RK or RK8.
- (3) Use RKF45 or GERK in a one step mode (this capability is permitted).

If option (3) is taken, then the user may just repeatedly call RKF45 or GERK (in the one step mode) until  $a_N$  is reached.

## ADAPTIVE ADAMS SOLUTION OF NONSTIFF DIFFERENTIAL EQUATIONS

Let  $y'(t) = f(t, y(t))$  denote a system of  $n$  ordinary first order differential equations where  $f(t, y) = (f_1(t, y), \dots, f_n(t, y))$  and  $y(t) = (y_1(t), \dots, y_n(t))$ . Assume that  $y(a)$  is known. Then for  $b \neq a$  the subroutine ODE is available for computing  $y(b)$ . ODE is recommended for nonstiff equations. The algorithm used is a variable order, variable step Adams predictor-corrector procedure.

**CALL ODE( $F, n, Y, T, TOUT, RERR, AERR, IND, WK, IWK$ )**

The argument  $F$  is the name of a user defined subroutine that has the format:

**CALL  $F(t, Y, DY)$**

$Y$  and  $DY$  are arrays of dimension  $n$ . On input  $Y$  contains the values  $y_1(t), \dots, y_n(t)$  for the argument  $t$ .  $F$  computes the derivatives  $y'_1(t), \dots, y'_n(t)$  using  $y'(t) = f(t, y(t))$  and stores the results in  $DY$ .  $F$  must be declared in the calling program to be of type EXTERNAL.

$WK$  is an array of dimension  $100 + 21n$  or larger, and  $IWK$  is an array of dimension 5 or larger.  $WK$  and  $IWK$  contain information needed for subsequent calls to ODE.

It is assumed that  $a \neq b$ . The argument  $Y$  of ODE is an array of dimension  $n$ , and the arguments  $T, RERR, AERR, IND$  are variables. ( $TOUT$  need not be a variable.) When ODE is initially called, it is assumed that:

$T = a$

$TOUT = b$

$Y(1), \dots, Y(n)$  contain the values  $y_1(a), \dots, y_n(a)$

$RERR$  = the relative error tolerance to be satisfied

$AERR$  = the absolute error tolerance to be satisfied

$IND = \pm 1$

It is preferable, both for efficiency and accuracy, that ODE be permitted to step along the axis from  $a$  to  $b$  using the largest steps that are appropriate. This is what is done when  $IND$  is set to 1. If  $IND = 1$  then ODE will step along the axis, possibly passing  $b$  and going as far as the point  $a + 10(b - a)$ . If  $b$  is passed, then the solution for the equations at  $b$  is obtained by interpolation. However,  $IND = 1$  cannot be used if the equations are not defined at all points between  $b$  and  $a + 10(b - a)$ . In a situation such as this, when integration cannot be permitted to step internally past  $TOUT$ ,  $IND$  must be set to  $-1$ . If  $IND = -1$  then it is required that the subroutine  $F$  be defined at  $TOUT$ . However,  $F$  need not be defined at points  $t$  past  $TOUT$ . If the equations  $y'(t) = f(t, y(t))$  are not defined at  $t = TOUT$ , then it should suffice to let  $F$  set each  $DY(i) = 0$  when  $t = TOUT$ . A solution (if one exists) will be obtained by extrapolation.

If  $IND$  is positive (negative), then when ODE terminates  $IND$  will have been reset by ODE to one of the values 2, 3, 4, 5, 6, 7 (2, -3, -4, -5, -6, 7). These values have the following meanings:

$IND = 2$  The equations have been solved at  $TOUT$ .  $T$  now has the value  $TOUT$  and  $Y$  contains the solution at  $TOUT$ .

$IND = \pm 3$  The error tolerances  $RERR$  and  $AERR$  are too small.  $T$  is set to the point closest to  $TOUT$  for which the equations were solved

and  $Y$  contains the solution at the point. RERR and AERR have been reset to larger acceptable values.

IND =  $\pm 4$  MAXNUM steps were performed.<sup>1</sup> More steps are needed to reach TOUT.  $T$  is set to the point closest to TOUT for which the equations were solved and  $Y$  contains the solution at the point.

IND =  $\pm 5$  MAXNUM steps were performed. More steps are needed to reach TOUT.  $T$  is set to the point closest to TOUT for which the equations were solved and  $Y$  contains the solution at the point. The equations appear to be stiff.

IND =  $\pm 6$  ODE did not reach TOUT because AERR = 0.  $T$  is set to the point closest to TOUT for which the equations were solved and  $Y$  contains the solution at the point.

IND = 7 No computation was performed. An input error was detected. The user must correct the error and call ODE again.

If IND =  $\pm 3, \pm 4, \pm 5$  then to continue the integration just call ODE again. Similarly, if IND =  $\pm 6$  then reset AERR to be positive and call ODE again. In these cases do not modify  $T, Y, \text{IND}$ . The output values for these parameters are the appropriate input values for the next call to ODE. However, AERR and RERR may always be modified when continuing an integration.

If the equations appear to be stiff (i.e., if IND =  $\pm 5$ ) then ODE may not be suitable for solving the equations. In this case it is recommended that a routine designed specifically for stiff equations be used.

Whenever IND = 2 occurs, then the equations have been solved at TOUT =  $b$ . WK and IWK contain information that can often be reused in continuing along the axis and solving the equations at a new point  $c$ . To continue the integration, normally one need only reset TOUT to the new value  $c$  and call ODE again. Do not modify  $T, Y, \text{IND}$ . The output values for these parameters are normally the appropriate input values for the next call to ODE. The one exception is when the equations are not defined at points past  $c$ . If this occurs, then one should also reset the output value IND = 2 (from the last call to ODE) to the input value IND = -2 for the next call to ODE. If IND is reset to -2, then integration will not proceed internally past the new TOUT when ODE is recalled. In this case, the subroutine  $F$  need not be defined for points past TOUT. However, it is required that  $F$  be defined at TOUT.

If after going from  $a$  to  $b$ , ODE is recalled to continue the integration and solve the equations at a new point  $c$ , then it is important that IND be set to  $\pm 2$  for the next call to ODE. Setting IND to  $\pm 1$  causes the integration procedure to be restarted, thereby eliminating the information being saved in WK and IWK. Restarting not only can take more time, but also can lead to less accurate results. If IND is set to  $\pm 2$ , then the integration procedure restarts itself only if the direction of integration is being reversed or IND was negative when ODE was last recalled. The direction of integration is reversed when  $b$  does not lie between  $a$  and  $c$ .

---

<sup>1</sup> Each step normally requires two calls to the subroutine  $F$ . Currently the internal parameter MAXNUM is set at 500.



If one has a choice between setting IND to be positive or negative, then always set IND to be positive. Extrapolation is normally involved when IND is negative. The extrapolation can require more time and be less accurate than the procedures employed when IND is positive.

**Input Errors.** IND = 7 when one of the following conditions is violated:

- (1)  $n \geq 1$
- (2)  $T \neq \text{TOUT}$
- (3)  $\text{RERR} \geq 0$  and  $\text{AERR} \geq 0$
- (4)  $\text{RERR}$  and  $\text{AERR}$  are not both 0
- (5)  $1 \leq |\text{IND}| \leq 5$ , or  $\text{IND} = \pm 6$  and  $\text{AERR} > 0$
- (6) When continuing an integration, the input value for  $T$  is the output value of TOUT from the previous call to ODE.

The last condition is automatically satisfied if the user has not inadvertently modified  $T$ .

**Error Control.** Assuming that ODE has obtained the correct value for  $y(t)$ , let  $e_i$  denote the error generated in the computation  $Y(i)$  of  $y_i(t+h)$  for  $i = 1, \dots, n$  when ODE steps from  $t$  to  $t+h$ . The routine attempts at each step to maintain the accuracy  $\sum_i (e_i/w_i)^2 \leq 1$  where  $w_i = \text{RERR} |Y(i)| + \text{AERR}$ . When this criterion is satisfied, then  $|e_i| \leq w_i$  for  $i = 1, \dots, n$ . This criterion includes as special cases relative error ( $\text{AERR} = 0$ ) and absolute error ( $\text{RERR} = 0$ ). However, if  $\text{AERR} = 0$  and  $Y(i) = 0$  for some  $i$ , then  $w_i = 0$  and  $\text{IND} = \pm 6$ .

When going from  $T$  to TOUT, ODE continually adjusts and readjusts its order and step size so as to maintain accuracy at each step. However, no attempt is made to control the progressive erosion of accuracy that can occur when the steps accumulate. Since the erosion of accuracy can be significant, at times one may wish to double-check the results by rerunning the problem. If this is done, then in the second run ask for greater accuracy.

**Programming.** ODE employs the subroutines DE1, STEP1, and INTRP. These routines were written by L. F. Shampine and M. K. Gordon (Sandia Laboratories). The function SPMPAR is also used.

**Reference.** Shampine, L. F., and Gordon, M. K., *Computer Solution of Ordinary Differential Equations*, W. H. Freeman and Company, San Francisco, 1975.

## ADAPTIVE RKF SOLUTION OF NONSTIFF DIFFERENTIAL EQUATIONS

Let  $y'(t) = f(t, y(t))$  denote a system of  $n$  ordinary first order differential equations where  $f(t, y) = (f_1(t, y), \dots, f_n(t, y))$  and  $y(t) = (y_1(t), \dots, y_n(t))$ . Assume that  $y(a)$  is known. Then for  $b \neq a$  the subroutine RKF45 is available for computing  $y(b)$ . RKF45 was designed for solving nonstiff differential equations when derivative evaluations are inexpensive and the accuracy requirements are low (less than 8 significant digits). The routine employs the fourth-fifth order Runge-Kutta-Fehlberg formulae.

**CALL RKF45(F, n, Y, T, TOUT, RERR, AERR, IND, WK, IWK)**

The argument  $F$  is the name of a user defined subroutine that has the format:

**CALL F(t, Y, DY)**

$Y$  and  $DY$  are arrays of dimension  $n$ . On input  $Y$  contains the values  $y_1(t), \dots, y_n(t)$  for the argument  $t$ .  $F$  computes the derivatives  $y'_1(t), \dots, y'_n(t)$  using  $y'(t) = f(t, y(t))$  and stores the results in  $DY$ .  $F$  must be declared in the calling program to be of type EXTERNAL.

$WK$  is an array of dimension  $3 + 6n$  or larger, and  $IWK$  is an array of dimension 5 or larger.  $WK$  and  $IWK$  contain information needed for subsequent calls to RKF45.

The argument  $Y$  of RKF45 is an array of dimension  $n$ , and the arguments  $T$ ,  $RERR$ ,  $IND$  are variables. ( $TOUT$  and  $AERR$  need not be variables.) When RKF45 is initially called, it is assumed that:

$T = a$

$TOUT = b$

$Y(1), \dots, Y(n)$  contain the values  $y_1(a), \dots, y_n(a)$

$RERR$  = the relative error tolerance to be satisfied

$AERR$  = the absolute error tolerance to be satisfied

$IND = \pm 1$

Normally  $IND = 1$ . However, if only a single step in the direction of  $TOUT$  is to be taken, then set  $IND = -1$ .

On output  $T$  is set to the point closest to  $TOUT$  for which the equations were solved, and  $Y$  contains the solution at  $T$ . Also  $IND$  reports the status of the results. RKF45 sets  $IND$  to one of the following values:

$IND = 2$  The equations were successfully solved at  $TOUT$ .  $T$  now has the value  $TOUT$ .

$IND = -2$  A single step in the direction of  $TOUT$  was taken.

$IND = 3$  The error tolerance  $RERR$  was too small.  $RERR$  has been reset to a larger acceptable value.

$IND = 4$  3000 derivative evaluations were performed. More derivative evaluations are needed to reach  $TOUT$ .

$IND = 5$  RKF45 did not reach  $TOUT$  because  $AERR = 0$ .  $AERR$  must be made positive.

$IND = 6$  Too much accuracy has been requested.  $AERR$  and/or  $RERR$  must be increased in value.

IND = 7 The closeness of the output points is restricting the natural step size choice.

IND = 8 No computation was performed. An input error was detected. The user must correct the error and call RKF45 again.

If IND = 2 then the equations have been solved at TOUT =  $b$ . The arrays WK and IWK contain information that can often be reused in continuing along the axis and solving the equations at a new point  $c$ . To continue the integration the user need only reset TOUT to the new point  $c$  and call RKF45 again.

If IND = -2 then to continue the integration another single step just call RKF45 again. In the single step mode (IND = -1, -2) the user must keep in mind that each step taken is in the direction of the current TOUT. Upon reaching TOUT (which is indicated by IND being set to 2), the user may then define a new TOUT and set IND to  $\pm 2$  for further integration.

If IND = 3 or 4 then to continue the integration just call RKF45 again. However, if IND = 5 then the user must first reset AERR to be positive before RKF45 can be recalled. If IND = 6 then it is required that IND be reset to  $\pm 2$  and that AERR and/or RERR be increased in value. If this is not done then the run will be terminated by a STOP instruction!

If IND = 7 then the user should either switch to another routine, or he should use the one step mode, setting IND = -2 for the next call to RKF45. This situation is discussed in the **Initial Value Solvers - Introductory Comments** section. If the user insists on continuing the integration with RKF45 in the standard multistep mode, then it is required that IND be reset to 2 before RKF45 is recalled. If this is not done then the run will be terminated by a STOP instruction.

If after going from  $a$  to  $b$ , RKF45 is recalled to continue the integration and solve the equations at a new point  $c$ , then it is important that IND be set to  $\pm 2$  instead of  $\pm 1$ . Setting IND =  $\pm 1$  causes the integration process to be restarted, thereby eliminating the information being saved in WK and IWK. Restarting wastes time and is normally not needed. The one exception is when the direction of integration is to be reversed. Then the integration must be restarted.

#### Notes.

- (1) AERR and RERR can be modified each time that RKF45 is called.
- (2) When continuing an integration, one may switch from the standard multistep mode (IND = 2) to the one step mode (IND = -2) whenever it is convenient to do so.

**Input Errors.** IND = 8 occurs when one of the following conditions is violated:

- (1)  $n \geq 1$
- (2)  $T \neq \text{TOUT}$  when  $\text{IND} \neq \pm 1$
- (3)  $\text{RERR} \geq 0$  and  $\text{AERR} \geq 0$
- (4)  $\text{IND} = \pm 1, \pm 2, 3, 4, \dots, 8$

**Error Control.** Pure absolute error control is not permitted. If RERR = 0 then RERR is reset to the smallest tolerance that is permitted for the computer being used, IND is set to 3, and the routine terminates.

When going from  $T$  to TOUT, RKF45 continually adjusts and readjusts its step size so as to maintain accuracy at each step. Assuming that RKF45 has obtained the correct value for  $y(t)$ , let  $e_i$  denote the error generated in the computation of  $y_i(t+h)$  for  $i = 1, \dots, n$  when RKF45 steps from  $t$  to  $t+h$ . Then at each step the error is controlled so that

$$|e_i| \leq \frac{|y_i(t)| + |y_i(t+h)|}{2} \text{RERR} + \text{AERR}$$

for  $i = 1, \dots, n$ . However, no attempt is made to control the progressive erosion of accuracy that can occur when the steps accumulate. Since the erosion of accuracy can be significant, at times one may wish to double-check the results. This can best be done by comparing the results obtained by RKF45 with those obtained by ODE or GERK. If ODE is used then ask for greater accuracy. However, if GERK is used then the current error tolerances can be used. GERK is more accurate than RKF45, and it estimates the global error generated.

**Programming.** RKF45 employs the subroutines RKFS and FEHL. These routines were written by H. A. Watts and L. F. Shampine (Sandia Laboratories). The function SPMPAR is also used.

**References.** Shampine, L. F., and Allen, R. C., *Numerical Computing: An Introduction*, W. B. Sanders, Philadelphia, 1973.

## ADAPTIVE RKF SOLUTION OF NONSTIFF DIFFERENTIAL EQUATIONS WITH GLOBAL ERROR ESTIMATION

Let  $y'(t) = f(t, y(t))$  denote a system of  $n$  ordinary first order differential equations where  $f(t, y) = (f_1(t, y), \dots, f_n(t, y))$  and  $y(t) = (y_1(t), \dots, y_n(t))$ . Assume that  $y(a)$  is known. Then for  $b \neq a$  the subroutine GERK is available for computing  $y(b)$ . GERK was designed for solving nonstiff differential equations when derivative evaluations are inexpensive and the accuracy requirements are low (less than 8 significant digits). The routine employs the fourth-fifth order Runge-Kutta-Fehlberg formulae. GERK estimates the accuracy of the solution  $y(b)$ .

**CALL GERK(*F, n, Y, T, TOUT, RERR, AERR, IND, GERROR, WK, IWK*)**

The argument  $F$  is the name of a user defined subroutine that has the format:

**CALL  $F(t, Y, DY)$**

$Y$  and  $DY$  are arrays of dimension  $n$ . On input  $Y$  contains the values  $y_1(t), \dots, y_n(t)$  for the argument  $t$ .  $F$  computes the derivatives  $y'_1(t), \dots, y'_n(t)$  using  $y'(t) = f(t, y(t))$  and stores the results in  $DY$ .  $F$  must be declared in the calling program to be of type EXTERNAL.

$WK$  is an array of dimension  $3 + 8n$  or larger, and  $IWK$  is an array of dimension 5 or larger.  $WK$  and  $IWK$  contain information needed for subsequent calls to GERK.

The argument  $Y$  of GERK is an array of dimension  $n$  or larger, and the arguments  $T$  and  $IND$  are variables. ( $TOUT$ ,  $RERR$ ,  $AERR$  need not be variables.) When GERK is initially called, it is assumed that:

$T = a$

$TOUT = b$

$Y(1), \dots, Y(n)$  contain the values  $y_1(a), \dots, y_n(a)$

$RERR$  = the relative error tolerance to be satisfied

$AERR$  = the absolute error tolerance to be satisfied

$IND = \pm 1$

Normally  $IND = 1$ . However, if only a single step in the direction of  $TOUT$  is to be taken, then set  $IND = -1$ .

$GERROR$  is an array of dimension  $n$  or larger. On output  $T$  is set to the point closest to  $TOUT$  for which the equations were solved,  $Y$  contains the solution at  $T$ , and  $GERROR(i)$  is an estimate of the error of  $Y(i)$  for  $i = 1, \dots, n$ . Also  $IND$  reports the status of the results. GERK sets  $IND$  to one of the following values:

$IND = 2$  The equations were successfully solved at  $TOUT$ .  $T$  now has the value  $TOUT$ .

$IND = -2$  A single step in the direction of  $TOUT$  was taken.

$IND = 3$  9000 derivative evaluations were performed. More derivative evaluations are needed to reach  $TOUT$ .

$IND = 4$  GERK did not reach  $TOUT$  because  $AERR = 0$ .  $AERR$  must be made positive.

- IND = 5 Too much accuracy has been requested. AERR and/or RERR must be increased in value.
- IND = 6 The closeness of the output points is restricting the natural step size choice.
- IND = 7 No computation was performed. An input error was detected. The user must correct the error and call GERK again.

If IND = 2 then the equations have been solved at  $TOUT = b$ . The arrays WK and IWK contain information that can often be reused in continuing along the axis and solving the equations at a new point  $c$ . To continue the integration the user need only reset TOUT to the new point  $c$  and call GERK again.

If IND = -2 then to continue the integration another single step just call GERK again. In the single step mode (IND = -1, -2) the user must keep in mind that each step taken is in the direction of the current TOUT. Upon reaching TOUT (which is indicated by IND being set to 2), the user may then define a new TOUT and set IND to  $\pm 2$  for further integration.

If IND = 3 then to continue the integration just call GERK again. However, if IND = 4 then the user must first reset AERR to be positive before GERK can be recalled. If IND = 5 then it is required that IND be reset to  $\pm 2$  and that AERR and/or RERR be increased in value. If this is not done then the run will be terminated by a STOP instruction!

If IND = 6 then the user should either switch to another routine, or he should use the one step mode, setting IND = -2 for the next call to GERK. This situation is discussed in the **Initial Value Solvers - Introductory Comments** section. If the user insists on continuing the integration with GERK in the standard multistep mode, then it is required that IND be reset to 2 before GERK is recalled. If this is not done then the run will be terminated by a STOP instruction.

If after going from  $a$  to  $b$ , GERK is recalled to continue the integration and solve the equations at a new point  $c$ , then it is important that IND be set to  $\pm 2$  instead of  $\pm 1$ . Setting IND =  $\pm 1$  causes the integration process to be restarted, thereby eliminating the information being saved in WK and IWK. Restarting wastes time and is normally not needed. The one exception is when the direction of integration is to be reversed. Then the integration must be restarted.

#### Notes.

- (1) AERR and RERR can be modified each time that GERK is called.
- (2) When continuing an integration, one may switch from the standard multistep mode (IND = 2) to the one step mode (IND = -2) whenever it is convenient to do so.

**Input Errors.** IND = 7 occurs when one of the following conditions is violated:

- (1)  $n \geq 1$
- (2)  $T \neq TOUT$  when  $IND \neq \pm 1$
- (3)  $RERR \geq 0$  and  $AERR \geq 0$
- (4)  $IND = \pm 1, \pm 2, 3, 4, \dots, 7$

**Accuracy Considerations.** Error control in GERK is almost identical to that in RKF45. One minor difference is that GERK never employs relative error tolerances less than  $3 \cdot 10^{-11}$ , whereas RKF45 never employs relative error tolerances less than  $10^{-12}$ .

The only significant difference between GERK and RKF45 is that GERK generates two solutions for the differential equations, whereas RKF45 generates only one. Let  $y(t)$  and  $\tilde{y}(t)$  denote the solutions generated by GERK at point  $t$ . One of these solutions, say  $y(t)$ , will frequently be identical to the solution computed by RKF45. When going from  $t$  to  $t + h$ , the step size  $h$  is selected so that  $y(t + h)$  satisfies the local error criterion. After a suitable  $h$  is found then GERK takes two steps, each of length  $h/2$ , to generate  $\tilde{y}(t + h)$  from  $\tilde{y}(t)$ . When GERK terminates, say at point  $T$ , then the  $\tilde{y}(T)$  solution is stored in the  $Y$  array and GERK estimates the error of  $\tilde{y}_i(T)$  to be  $(y_i(T) - \tilde{y}_i(T))/31$  for  $i = 1, \dots, n$ .

**Programming.** GERK employs the subroutines GERKS and FEHL. These routines were written by H. A. Watts and L. F. Shampine (Sandia Laboratories). The function SPMPAR is also used.

**Reference.** Shampine, L. F., and Allen, R. C., *Numerical Computing: An Introduction*, W. B. Saunders, Philadelphia, 1973.

## ADAPTIVE SOLUTION OF STIFF DIFFERENTIAL EQUATIONS

Let  $y'(t) = f(t, y(t))$  denote a system of  $n$  ordinary first order differential equations where  $f(t, y) = (f_1(t, y), \dots, f_n(t, y))$  and  $y(t) = (y_1(t), \dots, y_n(t))$ . Assume that  $y(a)$  is known. Then for  $b \neq a$  the following subroutines are available for computing  $y(b)$ . These routines are designed for stiff differential equations. The algorithm used is a variable order, variable step backward differentiation procedure.

**CALL SFODE**( $F, n, Y, T, TOUT, INFO, RERR, AERR, IER,$   
 $WK, \ell, IWK, m, RD, ID$ )

$RD$  and  $ID$  are arrays defined by the user containing any real and integer data that is needed for computing  $f$ . These arrays may contain any information that the user desires. The argument  $F$  is the name of a user defined subroutine that has the format:

**CALL**  $F(t, Y, DY, RD, ID)$

$Y$  and  $DY$  are arrays of dimension  $n$ . On input  $Y$  contains the values  $y_1(t), \dots, y_n(t)$  for the argument  $t$ .  $F$  computes the derivatives  $y'_1(t), \dots, y'_n(t)$  using  $y'(t) = f(t, y(t))$  and stores the results in  $DY$ .  $F$  must be declared in the calling program to be of type EXTERNAL.

$INFO$  is an array of dimension 4,  $WK$  an array of dimension  $\ell$ , and  $IWK$  an array of dimension  $m$ .  $WK$  and  $IWK$  are work spaces for the routine, and  $INFO$  is an array defined by the user containing information on how the equations are to be treated.

**INFO(1):** Set  $INFO(1) = 0$  on an initial call to the routine. On a continuation call  $INFO(1) = 1$ .

**INFO(2):** Normally  $INFO(2) = 0$ . However,  $INFO(2) = 1$  when the intermediate output mode is desired (see below).

**INFO(3):** When  $INFO(3) = 0$ , SFODE proceeds from  $a$  to  $b$  using the largest steps that are appropriate. If  $b$  is passed then  $y(b)$  is obtained by interpolation. However, for some problems the routine cannot be permitted to step past a point  $TSTOP$  because  $y'(t) = f(t, y(t))$  is discontinuous or not defined beyond  $TSTOP$ . When this is the case set  $INFO(3) = 1$  and  $WK(1) = TSTOP$ .

**INFO(4):** When proceeding from  $a$  to  $b$ , the  $n \times n$  Jacobian matrix  $Jf(t) = (\partial f_i / \partial y_j)$  is computed and stored in  $WK$ . Normally it is assumed that

$$\begin{aligned} INFO(4) &= 0 \\ \ell &\geq 250 + 10n + n^2 \\ m &\geq 55 + n. \end{aligned}$$

However, if  $Jf(t)$  is banded for all  $t$ , having the lower and upper band widths  $m_\ell$  and  $m_u$  where  $2m_\ell + m_u < n$ , then the following setup can be used:

$$\begin{aligned} INFO(4) &= 1 \\ IWK(1) &= m_\ell \\ IWK(2) &= m_u \\ \ell &\geq 250 + 10n + (2m_\ell + m_u + 1)n \\ m &\geq 55 + n \end{aligned}$$



$T$ ,  $TOUT$ ,  $RERR$ , and  $AERR$  are variables, and the argument  $Y$  of SFODE is an array of dimension  $n$ . On an initial call to the routine it is assumed that

$INFO(1) = 0$

$T = a$

$TOUT = b$

$Y(1), \dots, Y(n)$  contain the values  $y_1(a), \dots, y_n(a)$

$RERR$  = the relative error tolerance to be satisfied ( $RERR \geq 0$ )

$AERR$  = the absolute error tolerance to be satisfied ( $AERR \geq 0$ ).

$IER$  is a variable. When SFODE terminates  $T$  is the final point where the equations were solved,  $Y$  contains the solution at  $T$ , and  $IER$  reports the status of the results.  $IER$  is assigned one of the following values:

$IER = 1$  A step was taken in the intermediate output mode.  $TOUT$  was not reached. To continue, call the routine again.

$IER = 2$  The solution at  $TOUT$  was obtained by stepping exactly to  $TOUT$ .

$IER = 3$  The solution at  $TOUT$  was obtained by stepping past  $TOUT$  and then interpolating. On output  $T = TOUT$ .

$IER = -1$  500 steps have been taken.  $TOUT$  has not been reached. To continue, call the routine again.

$IER = -2$  The tolerances  $RERR$  and  $AERR$  were too stringent.  $RERR$  and  $AERR$  have been modified by the routine. The tolerances may be further modified by the user if he desires. To continue, call the routine again.

$IER = -3$  In this case  $AERR = 0$ . SFODE stopped when  $y_i$  became 0.  $INFO(1)$  was set to  $-i$ . To continue set  $AERR$  to be positive,  $INFO(1) = 1$ , and call the routine again.

$IER = -6$  Convergence failed on the last attempted step. An inaccurate Jacobian matrix may be the problem. To continue, restart the routine by setting  $INFO(1) = 0$  and call the routine again.

$IER = -7$  Repeated error test failures occurred on the last attempted step. The problem should be reexamined. A singularity may be present in the solution. To continue, restart by setting  $INFO(1) = 0$  and call the routine again.

$IER \leq -33$  An input error was detected (see below).

When  $IER \geq -2$ , then  $INFO(1) = 1$  on output.

When the equations are solved at  $TOUT$  ( $IER = 2$  or  $3$ ), integration can be continued along the axis to solve the equations at a new point  $c$  beyond  $TOUT$ . To continue, one need only set  $TOUT$  to the new value  $c$  and call the routine again. When continuing an integration where  $INFO(1) = 1$ , never modify  $T$ ,  $Y$ ,  $WK$ ,  $IWK$ ,  $INFO(3)$ , and  $INFO(4)$ . However,  $INFO(2)$ ,  $RERR$ ,  $AERR$ ,  $RD$ , and  $ID$  may be modified at any time.

**Intermediate Output Mode.** If one wishes to study the behavior of the solution  $y(t)$  as the routine steps from  $T$  to  $TOUT$ , then set  $INFO(2) = 1$ . Then SFODE will stop after each successful step (yielding  $IER = 1$ ) until  $TOUT$  is reached. One may switch from the standard mode of operation ( $INFO(2) = 0$ ) to the intermediate output mode ( $INFO(2) = 1$ ) or visa versa at any time.

**Remark.** The diagnostic  $IER = -1$  does not state that 500 steps have been taken on the current call to SFODE. On an initial call to the routine the step counter is set to 0. On continuation calls, the counter continues to increase until 500 steps have accumulated. When  $IER = -1$  is reported, the counter is reset to 0, and only then does the step counting begin again.

**Input Errors.** IER is set to one of the following values when an input error is detected.

IER = -33  $n < 1$

IER = -34  $RERR < 0$

IER = -35  $AERR < 0$

IER = -36 The routine has been called with TOUT, but it has also been told not to step past the point TSTOP.

IER = -37  $T = TOUT$ . This is not permitted on continuation calls.

IER = -38 The user has modified  $T$ .

IER = -39 TOUT is not beyond  $T$ . An attempt is being made to change the direction of integration without restarting.

IER = -40 The Jacobian matrix is banded. However,  $m_l$  and  $m_u$  do not satisfy  $0 \leq m_l < n$  and  $0 \leq m_u < n$ .

IER = -41  $\ell < 250 + 10n + n^2$

IER = -42  $\ell < 250 + 10n + (2m_l + m_u + 1)n$

IER = -43  $m < 55 + n$

IER = -44 INFO(1) is incorrect.

After the error is corrected, set  $INFO(1) = 0$  and call the routine again.

**Error Control.** Assuming that SFODE has the correct value for  $y(t)$ , let  $e_i$  denote the error generated in computing  $y_i(t+h)$  for  $i = 1, \dots, n$  when SFODE steps from  $t$  to  $t+h$ . The routine attempts at each step to maintain the accuracy  $\frac{1}{n} \sum_i (e_i/w_i)^2 \leq 1$  where  $w_i = RERR |y_i(t)| + AERR$ . When this criterion is satisfied,  $|e_i| \leq \sqrt{n} w_i$  for  $i = 1, \dots, n$ . This criterion includes as special cases relative error ( $AERR = 0$ ) and absolute error ( $RERR = 0$ ). However, if  $AERR = 0$  and  $y_i(t) = 0$  for some  $i$ , then this criterion cannot be applied and  $IER = -3$  occurs.

When proceeding from  $T$  to TOUT, the routine continually readjusts its order and step size so as to maintain accuracy at each step. However, no attempt is made to control the progressive erosion of accuracy that can occur when the steps accumulate. Since the erosion of accuracy can be significant, at times one may wish to double-check the results. If the problem is nonstiff or mildly stiff for an interval, then the best procedure is to compare the results obtained by SFODE with those obtained by ODE for the interval. ODE normally maintains greater accuracy than SFODE. However, if the problem is extremely stiff then rerun the problem with SFODE. On the second run, request greater accuracy.

**Programming.** SFODE calls the subroutines STFODE and ZZZJAC. STFODE employs the subroutines LSOD1, HSTART, INTYD, STOD, CFOD, PJAC, SLVS, SGBFA, SGBSL, SGEFA, SGESL, SAXPY, and SSCAL, and the functions VNORM, VNWRMS, ISAMAX, SDOT, and SPMPAR. The routines save and exchange information in a labeled common block having the block name DEBDF1. STFODE is a modification by A. H. Morris of the

subroutine DEBDF, designed by L. F. Shampine and H. A. Watts (Sandia Laboratories). DEBDF appears in the SLATEC library. SFODE is a driver for a modification of the code LSODE, written by A. C. Hindmarsh (Lawrence Livermore Laboratory).

continuous over  $t$

**CALL SFODE1**( $F, n, Y, T, TOUT, INFO, RERR, AERR, IER,$   
 $WK, \ell, IWK, m, RD, ID$ )

SFODE1 differs from SFODE only in the treatment of RERR and AERR. In SFODE1, RERR and AERR are arrays of dimension  $n$ . RERR( $i$ ) and AERR( $i$ ) are relative and absolute error tolerances to control the accuracy of the  $i^{th}$  solution component  $y_i(t)$  for  $i = 1, \dots, n$ . Let  $e_i$  denote the error generated in the computation of  $y_i(t+h)$  from  $y_i(t)$  when SFODE1 steps from  $t$  to  $t+h$ . Then SFODE1 attempts at each step to maintain the accuracy  $\frac{1}{n} \sum_i (e_i/w_i)^2 \leq 1$  where  $w_i = RERR(i)|y_i(t)| + AERR(i)$ . When this criterion is satisfied  $|e_i| \leq \sqrt{n}w_i$  for  $i = 1, \dots, n$ . However, if AERR( $i$ ) = 0 and  $y_i(t) = 0$  for some  $i$ , then the criterion cannot be applied and IER = -3 occurs.

When IER references RERR and AERR, the settings for IER provide the following information:

IER = -2 The accuracy required by RERR and AERR is too stringent. RERR and AERR have been modified by the routine. RERR and AERR may be further modified by the user if he desires. To continue, call SFODE1 again.

IER = -3 SFODE1 stopped when  $y_i$  became 0 and AERR( $i$ ) = 0. INFO(1) was set to  $-i$ . To continue set AERR( $i$ ) to be positive, INFO(1) = 1, and call the routine again.

IER = -34 (Input error) RERR( $i$ ) < 0 for some  $i$ .

IER = -35 (Input error) AERR( $i$ ) < 0 for some  $i$ .

RERR and AERR may be modified on any continuation call to SFODE1.

**Programming.** SFODE1 calls the subroutines STFODE and ZZZJAC.

## FOURTH-ORDER RUNGE-KUTTA

Let  $y'(t) = f(t, y(t))$  denote a system of  $n$  ordinary first order differential equations where  $f(t, y) = (f_1(t, y), \dots, f_n(t, y))$  and  $y(t) = (y_1(t), \dots, y_n(t))$ . Assume that  $y(t_0)$  is known. Then for a small real number  $h$ , the subroutine RK is available for computing  $y(t_0 + h)$ . RK employs the standard fourth-order Runge-Kutta procedure.

**CALL RK( $n, T, h, A, F$ )**

The argument  $F$  is the name of a user defined subroutine that has the format:

**CALL F( $t, Z$ )**

$Z$  is an array of dimension  $n$  containing the values  $y_1(t), \dots, y_n(t)$  for the argument  $t$ .  $F$  computes the derivatives  $y'_1(t), \dots, y'_n(t)$  using  $y'(t) = f(t, y(t))$  and stores the results in  $Z$ , thereby destroying the original data in  $Z$ .  $F$  must be declared in the calling program to be of type EXTERNAL.

$T$  is a variable having the value  $t_0$  and  $A$  an array of dimension  $3n$  or larger. It is assumed that  $A(1), \dots, A(n)$  contain the values  $y_1(t_0), \dots, y_n(t_0)$ . If  $h = 0$  then RK computes the derivatives  $y'_1(t_0), \dots, y'_n(t_0)$  and stores them in  $A(n+1), \dots, A(2n)$ . If  $h \neq 0$  then it is assumed that the derivatives  $y'_1(t_0), \dots, y'_n(t_0)$  have already been computed and stored in  $A(n+1), \dots, A(2n)$ . In this case, when RK is called, the values  $y_1(t_0 + h), \dots, y_n(t_0 + h)$  and derivatives  $y'_1(t_0 + h), \dots, y'_n(t_0 + h)$  are computed and stored in  $A(1), \dots, A(2n)$ . Also  $T$  is reset to the value  $t_0 + h$ .

**Note.** The area  $A(2n+1), \dots, A(3n)$  serves as work space for the routine.

**Example.** Consider the equations

$$x'(t) = y(t)$$

$$y'(t) = -x(t)$$

where  $x(0) = 0$  and  $y(0) = 1$ . The following code may be used for solving these equations at the points .01, .02, ..., 1.00, and storing the results in the arrays  $X$  and  $Y$ .

```

DIMENSION A(6), X(100), Y(100)
EXTERNAL FUN
T = 0.0
H = .01
A(1) = 0.0
A(2) = 1.0
A(3) = 1.0
A(4) = 0.0
DO 10 I = 1, 100
CALL RK(2,T,H,A,FUN)
X(I) = A(1)
10 Y(I) = A(2)

```

Here FUN may be defined by:

```
SUBROUTINE FUN(T,Z)
  DIMENSION Z(2)
  X = Z(1)
  Y = Z(2)
  Z(1) = Y
  Z(2) = -X
  RETURN
END
```

Note that the statements  $A(3) = 1.0$  and  $A(4) = 0.0$ , which store the derivatives  $x'(0)$  and  $y'(0)$  in  $A(3)$  and  $A(4)$ , can be replaced with `CALL RK(2,T,0.0,A,FUN)`.

**Programmer.** A.H. Morris.

## EIGHTH-ORDER RUNGE-KUTTA

Let  $y'(t) = f(t, y(t))$  denote a system of  $n$  ordinary first order differential equations where  $f(t, y) = (f_1(t, y), \dots, f_n(t, y))$  and  $y(t) = (y_1(t), \dots, y_n(t))$ . Assume that  $y(t_0)$  is known. Then for a small real number  $h$ , the subroutine RK8 is available for computing  $y(t_0 + h)$ .

**CALL RK8**( $n, T, h, Y, DY, WK, F$ )

The argument  $F$  is the name of a user defined subroutine that has the format:

**CALL**  $F(t, Z)$

$Z$  is an array of dimension  $n$  containing the values  $y_1(t), \dots, y_n(t)$  for the argument  $t$ .  $F$  computes the derivatives  $y'_1(t), \dots, y'_n(t)$  using  $y'(t) = f(t, y(t))$  and stores the results in  $Z$ , thereby destroying the original data in  $Z$ .  $F$  must be declared in the calling program to be of type EXTERNAL.

$T$  is a variable having the value  $t_0$ , and  $Y$  and  $DY$  are arrays of dimension  $n$ . It is assumed that  $Y$  contains the values  $y_1(t_0), \dots, y_n(t_0)$ . If  $h = 0$  then RK8 computes the derivatives  $y'_1(t_0), \dots, y'_n(t_0)$  and stores them in  $DY$ . If  $h \neq 0$  then it is assumed that the derivatives  $y'_1(t_0), \dots, y'_n(t_0)$  have already been computed and stored in  $DY$ . In this case, when RK8 is called, the values  $y_1(t_0 + h), \dots, y_n(t_0 + h)$  and derivatives  $y'_1(t_0 + h), \dots, y'_n(t_0 + h)$  are computed and stored in  $Y$  and  $DY$ , thereby destroying the original data in  $Y$  and  $DY$ . Also  $T$  is reset to the value  $t_0 + h$ .

$WK$  is an array of dimension  $8n$  or larger that is used for a work space by the routine.

**Algorithm.** The routine employs formulae (8-12) given on p. 34 of the reference.

**Remarks.** RK8 is used in the same manner as the routine RK. RK8 takes more time and storage than RK, but may be more accurate.

**Reference.** Shanks, E. B., "Solutions of Differential Equations by Evaluations of Functions," *Math. Comp.* 20 (1966), pp. 21-38.

**Programmer.** A. H. Morris

## SEPARABLE SECOND-ORDER ELLIPTIC EQUATIONS ON RECTANGULAR DOMAINS

Given a separable elliptic equation

$$a(x)u_{xx} + b(x)u_x + c(x)u + d(y)u_{yy} + e(y)u_y + f(y)u = g(x, y)$$

on the rectangle  $a_1 \leq x \leq a_2, b_1 \leq y \leq b_2$ , where  $u$  is periodic in  $x$  or  $y$ , or  $u$  or its normal derivative  $\partial u / \partial n$  is given on each of the edges. For  $m, n > 1$  let  $x_i = a_1 + (i - 1)h$  and  $y_j = b_1 + (j - 1)k$  where  $h = (a_2 - a_1)/(m - 1), k = (b_2 - b_1)/(n - 1), i = 1, \dots, m$ , and  $j = 1, \dots, n$ . Then the subroutine SEPDE is available for computing  $u$  at the points  $(x_i, y_j)$ .

**CALL SEPDE(COFX,COFY,g,ITYPE,BVAL,IORD,a<sub>1</sub>,a<sub>2</sub>,m,b<sub>1</sub>,b<sub>2</sub>,n,  
U,ku,W,ℓ,IND)**

It is assumed that  $m \geq 7$  and  $n \geq 6$ .  $U$  is an  $m \times n$  matrix. The argument  $ku$  is the number of rows in the dimension statement for  $U$  in the calling program. When SEPDE is called, if the elliptic equation is solved then  $U(i, j) = u(x_i, y_j)$  for  $i = 1, \dots, m$  and  $j = 1, \dots, n$ .

The input argument IORD is the order of the approximation procedure to be used. IORD may have the values 2 or 4.

The argument COFX is the name of a user defined subroutine that has the format:

CALL COFX( $x, A, B, C$ )

$A, B$ , and  $C$  are variables. COFX sets  $A = a(x)$ ,  $B = b(x)$ , and  $C = c(x)$  for the argument  $x$ . COFX must be declared in the calling program to be of type EXTERNAL.

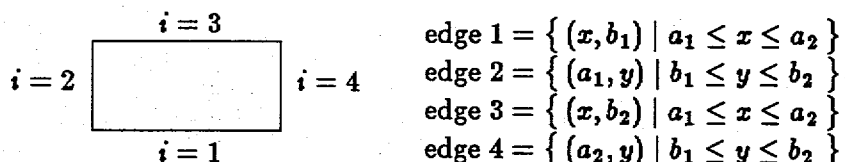
The argument COFY is the name of a user defined subroutine that has the format:

CALL COFY( $y, D, E, F$ )

$D, E$ , and  $F$  are variables. COFY sets  $D = d(y)$ ,  $E = e(y)$ , and  $F = f(y)$  for the argument  $y$ . COFY must be declared in the calling program to be of type EXTERNAL.

The argument  $g$  is the name of a user defined function, where  $g(x, y)$  gives the right hand side of the elliptic equation for all  $a_1 \leq x \leq a_2, b_1 \leq y \leq b_2$ . The argument  $g$  must be declared in the calling program to be of type EXTERNAL.

**Boundary Conditions.** The edges of the rectangular domain are labeled in a clockwise manner as follows:



ITYPE is an array of dimension 4. For edge  $i$  ( $i = 1, \dots, 4$ ), ITYPE( $i$ ) specifies the type of boundary condition on the edge. ITYPE( $i$ ) must be set by the user to one of the following values:

ITYPE( $i$ ) = 0 It is assumed that  $u$  is given on the edge.

ITYPE( $i$ ) = 1 If  $i = 1$  or  $i = 3$  then  $u_y$  is given on the edge. Otherwise, if  $i = 2$  or  $i = 4$  then  $u_x$  is given on the edge.

ITYPE( $i$ ) = -1 If  $i = 1$  or  $i = 3$  then it is assumed that  $u$  is periodic in  $y$ ; i.e.,  $u(x, y + b_2 - b_1) = u(x, y)$  for all  $x, y$ . In this case ITYPE( $i$ ) must be -1 for both  $i = 1$  and  $i = 3$ . If  $i = 2$  or  $i = 4$  then it is assumed that  $u$  is periodic in  $x$ ; i.e.,  $u(x + a_2 - a_1, y) = u(x, y)$  for all  $x, y$ . In this case ITYPE( $i$ ) must be -1 for both  $i = 2$  and  $i = 4$ .

The argument BVAL is the name of a user defined function. BVAL( $i, x, y$ ) is defined for any point  $(x, y)$  on edge  $i$  when ITYPE( $i$ ) = 0 or 1, where

$$\text{BVAL}(i, x, y) = \begin{cases} u(x, y) & \text{if ITYPE}(i) = 0 \\ u_y(x, y) & \text{if ITYPE}(i) = 1 \text{ (} i = 1 \text{ or } i = 3 \text{)} \\ u_x(x, y) & \text{if ITYPE}(i) = 1 \text{ (} i = 2 \text{ or } i = 4 \text{)} \end{cases}$$

The function BVAL( $i, x, y$ ) is ignored when ITYPE( $i$ ) = -1. BVAL must be declared in the calling program to be of type EXTERNAL.

$W$  is an array of dimension  $\ell$  that is a work space. The argument  $\ell$  is a variable whose value depends on IORD,  $m, n$ , and the types of boundary conditions used. Let  $\nu$  be the largest integer  $\leq \log_2 n$  and  $\ell_1 = (\nu - 1)2^{\nu+2} + \nu + 14m + 12n + 6$ . Then

$$\begin{aligned} \ell &\geq \ell_1 && \text{if IORD} = 2. \\ \ell &\geq \ell_1 + mn && \text{if IORD} = 4. \end{aligned}$$

When the routine terminates,  $\ell$  will have been reset to the actual amount of storage needed.

IND is a variable that reports the status of the results. When SEPDE terminates, IND has one of the following values:

- IND = 0 The solution  $U$  was obtained.
- IND = -1 A constant (which is stored in  $W(1)$ ) was subtracted from the right hand side of the equation in order to obtain a solution  $U$ . The solution is a weighted minimal least squares solution of the original problem.
- IND = 1 (Input error)  $a_1 \geq a_2$  or  $b_1 \geq b_2$ .
- IND = 2 (Input error) ITYPE( $i$ )  $\neq 0, \pm 1$  for some edge  $i$ .
- IND = 4 The approximating linear system of equations is not diagonally dominant. This cannot occur when  $m$  and  $n$  are sufficiently large. Increase  $m$  and  $n$ , and reset  $\ell$ .
- IND = 5 (Input error)  $ku < m$ .
- IND = 6 (Input error)  $m < 7$ .
- IND = 7 (Input error)  $n < 6$ .
- IND = 8 (Input error) IORD  $\neq 2, 4$ .
- IND = 10 (Input error)  $a(x)d(y) \leq 0$  for some interior point  $(x, y)$  of the rectangle. This violates the assumption that the equation is elliptic.
- IND = 11 (Input error)  $\ell$  was too small.  $\ell$  has been reset to the minimum amount of storage needed for  $W$ .



IND = 12 (Input error) ITYPE(*i*) = -1 for edge 1 or 3, but not for both edges.

IND = 13 (Input error) ITYPE(*i*) = -1 for edge 2 or 4, but not for both edges.

**Precision.** If IORD = 2 then the elliptic equation is approximated by a set of linear equations using finite differences. Otherwise, if IORD = 4 then the approximating equations are obtained by deferred corrections. The most accuracy is achieved when ITYPE(*i*) = 1 boundary conditions are not involved. For  $m, n \geq 100$ , 3-4 digit accuracy may be attained when IORD = 2 and 7-8 digit accuracy when IORD = 4. When ITYPE(*i*) = 1 boundary conditions are used, then for  $m, n \geq 100$ , 2-3 digits may be attained when IORD = 2 and 5-6 digits when IORD = 4.

**Programming.** SEPDE is an interface by A. H. Morris for SEPELL, a modification of the routine SEPELI described in the reference. SEPELI was developed by John C. Adams, being supported (in part) by codes written by Paul Swarztrauber and Roland Sweet (National Center for Atmospheric Research, Boulder, Colorado). SEPDE employs the subroutines PEDGE, SEPELL, SEPELI, CHKPRM, CHKSNG, ORTHG, MINSOL, TRISP, DEFER, DXFN, DYFN, BLKTRI, BLKTRI, COMPB, PROD0, PRODP, CPROD0, CPRODP, INDXA, INDXB, INDXC, PPADD, TQLRT0 and functions PSGF, BSRH, PPSGF, PPSPF, SPMPAR. The routines exchange information in the labeled common blocks having block names CBLKT and SPLP.

**Example.** Consider  $(1+x)^2 u_{xx} - 2(1+x)u_x + u_{yy} = 3(1+x)^4 \sin y$  for  $0 \leq x \leq 1$  and  $|y| \leq \pi$

where  $u_x(0, y) = 4 \sin y$   $|y| \leq \pi$   
 $u(1, y) = 16 \sin y$

and  $u$  is periodic in  $y$ . This problem has the solution  $u = (1+x)^4 \sin y$ . Let

ITYPE(1) = -1  
ITYPE(2) = 1  
ITYPE(3) = -1  
ITYPE(4) = 0.

Then the following routines and functions may be used for describing the problem. (Here  $g = \text{GVAL.}$ )

```
SUBROUTINE COFX (X,A,B,C)
  T = 1.0 + X
  A = T*T
  B = -2.0*T
  C = 0.0
  RETURN
END
```

```
SUBROUTINE COFY (Y,D,E,F)
```

```
D = 1.0  
E = 0.0  
F = 0.0  
RETURN  
END
```

```
REAL FUNCTION GVAL (X,Y)  
GVAL = 3.0*(1.0 + X)** 4*SIN(Y)  
RETURN  
END
```

```
REAL FUNCTION BVAL (I,X,Y)  
BVAL = 4.0*SIN(Y)  
IF (I .EQ. 4)BVAL = 4.0*BVAL  
RETURN  
END
```

COFX, COFY, GVAL, and BVAL must be declared in the calling program to be of type EXTERNAL.

**Reference.** Adams, J., Swarztrauber, P., and Sweet, R., *FISHPAK: Efficient FORTRAN Subprograms for the Solution of Separable Elliptic Partial Differential Equations, Version 3*. National Center for Atmospheric Research, Boulder, Colorado, 1978.

## UNIFORM RANDOM NUMBER GENERATOR

The following subroutine is available for generating a sequence of uniform variates in the interval  $(0, 1)$ .

**CALL URNG**(ix,A,n,IERR)

The argument  $n$  is the number of variates to be generated.  $A$  is an array of dimension  $n$  or larger, and  $ix$  and  $IERR$  are variables. On input,  $ix$  is an integer (called a *seed*) for initializing the sequence of variates. It is assumed that  $1 \leq ix < 2^{31} - 1$ . When URNG is called, if no input errors are detected then  $IERR$  is set to 0 and  $n$  uniform variates are stored in  $A$ . On output,  $ix$  is a new seed for generating more variates.

**Error Return.**  $IERR = 1$  if  $n \leq 0$  and  $IERR = 2$  if  $ix$  is not a proper seed.

**Usage.** A given seed always initiates the same set of variates. Thus, the following two sets of instructions

- (1)      $IX = 103$   
          CALL URNG( $IX, A, 30, IERR$ )
  
- (2)      $IX = 103$   
          CALL URNG( $IX, A, 20, IERR$ )  
          CALL URNG( $IX, A(21), 10, IERR$ )

generate the same 30 variates.

**Remark.** It is assumed that the integer arithmetic being used handles all integers  $i$  in the interval  $|i| \leq 2^{31} - 1$ .

**Programming.** Written by Linus Schrage (University of Chicago). Adapted by A.H. Morris.

**Reference.** Schrage, Linus, "A More Portable Fortran Random Number Generator," *ACM Trans. Math Software* 5 (1979), pp. 132-138.

## GAUSSIAN RANDOM NUMBER GENERATOR USING THE BOX-MULLER TRANSFORMATION

The following subroutine is available for generating a sequence of normal variates from a normal distribution with mean 0 and standard deviation 1.

**CALL NRNG**(ix,A,n,IERR)

The argument  $n$  is the number of variates to be generated.  $A$  is an array of dimension  $n$  or larger, and  $ix$  and  $IERR$  are variables. On input,  $ix$  is an integer (called a *seed*) for initializing the sequence of variates. It is assumed that  $1 \leq ix < 2^{31} - 1$ . When NRNG is called, if no input errors are detected then  $IERR$  is set to 0 and  $n$  normal variates are stored in  $A$ . On output,  $ix$  is a new seed for generating more variates.

**Error Return.**  $IERR = 1$  if  $n \leq 0$  and  $IERR = 2$  if  $ix$  is not a proper seed.

**Algorithm.** When NRNG is called, an even number of uniform variates  $u_1, \dots, u_m$  is generated ( $m = n$  if  $n$  is even and  $m = n + 1$  if  $n$  is odd). Then the Box-Muller transformation

$$\begin{aligned}a_k &= \sqrt{-2 \ln u_k} \cos 2\pi u_{k+1} \\a_{k+1} &= \sqrt{-2 \ln u_k} \sin 2\pi u_{k+1} \quad (k = 1, 3, 5, \dots)\end{aligned}$$

is applied to obtain  $n$  normal variates  $a_1, \dots, a_n$ . This transformation generates pairs of normal variates  $(a_1, a_2), (a_3, a_4), \dots$  from the corresponding pairs of uniform variates. If  $n$  is odd then only the first variate of the final pair  $(a_n, a_{n+1})$  is computed.

**Usage.** A given seed always generates the same sequence of uniform variates. Thus, if we consider the following three sets of instructions

- (1)      $IX = 73$   
          CALL NRNG( $IX, A, 30, IERR$ )
- (2)      $IX = 73$   
          CALL NRNG( $IX, A, 10, IERR$ )  
          CALL NRNG( $IX, A(11), 20, IERR$ )
- (3)      $IX = 73$   
          CALL NRNG( $IX, A, 9, IERR$ )  
          CALL NRNG( $IX, A(10), 20, IERR$ ),

then we note that (1) and (2) generate the same 30 normal variates. Also (3) produces 29 of these 30 variates, skipping the 10<sup>th</sup> normal variate. The reason for this is that the request in (3) for 9 normal variates requires 10 uniform variates to be generated. The 10 uniform variates could be used for computing 10 normal variates (as is done in (2)). However, since only 9 normal variates are requested, computation of the 10<sup>th</sup> normal variate that would usually be generated is skipped.

**Remark.** NRNG calls the subroutine URNG for the uniform variates. Thus, it is assumed that the integer arithmetic being used handles all integers  $i$  in the interval  $|i| \leq 2^{31} - 1$ .

**Programmer.** A.H. Morris

## APPENDIX

### Installation Of The NSWC Library

Two versions of the code for the NSWC library are maintained, differing only in the declaration of assumed size arrays. In the 1966 version, statements such as `REAL A(1)` are frequently used in declaring as arrays arguments *A* of a function or subroutine. In the 1977 version, statements such as `REAL A(*)` are used for declaring the arrays. The 1966 version must be used by all Fortran compilers which satisfy the 1966 standard, and it may be used by almost all the compilers which satisfy the 1977 standard. The 1977 version can only be used by the compilers which satisfy the 1977 standard.

Code for the library can be obtained on 7-track and 9-track tapes, and on  $5\frac{1}{4}$  inch disks that can be read by the IBM PC. Two files are given on a tape. The first file contains the 1966 version of the code, and the second file contains the 1977 version. The 1977 version is normally provided on the disks.

The first function in the NSWC library, namely `IPMPAR`, must be modified for the particular Fortran being used. `IPMPAR` provides the integer constants which characterize the integer, single precision, and double precision arithmetics being used (see pp. 3-4).

Instructions are given in the in-line documentation of `IPMPAR` for defining the constants that are needed. If constants are not provided for the compiler being used, then the Fortran manual for the compiler normally gives the constants for the integer arithmetic. However, help may be needed in obtaining the constants for the single and double precision arithmetics. The subroutines `MACH` and `RADIX` are provided for this purpose.

`MACH` and `RADIX`, and their subroutines `MACH1`, `STORE2`, `MACH2`, `DSTOR2` are the next subprograms after `IPMPAR`. Instructions for the use of `MACH` and `RADIX` are given in `MACH`. These subroutines are experimental. They are provided only as an aid for obtaining the constants for the single and double precision arithmetics. They are not used by any of the functions and subroutines in the NSWC library, and they are not considered to be part of the library. `MACH1` and `MACH2` perform some writing. None of the functions and subroutines in the NSWC library perform I/O.

After `IPMPAR` has been defined, the remainder of the library can be installed. None of the remaining code needs any modification. Codes for the functions and subroutines appear approximately in the order that they appear in the manual.

# INDEX

ABSLV ... 221	CDET ... 219	CSVDC ... 217	DERFC1 ... 39
ADAPT ... 377	CDOTC ... 151	CSWAP ... 147	DET ... 219
AI ... 84	CDOTU ... 151	CTIP ... 175	DGAMLN ... 63
AIE ... 84	CEIG ... 313	CTPOSE ... 175	DGAMMA ... 63
ALI ... 60	CEIGV ... 315	CTPRD ... 283	DHFTI ... 330
ALNREL ... 23	CERF ... 35	CTPRD1 ... 283	DHFTI2 ... 330
AORD ... 5	CERFC ... 36	CTTRANS ... 177	DISORT ... 6
ARCECO ... 199	CEXPLI ... 53	CTSLV ... 298	DKPROD ... 191
ARCESL ... 200	CFRNLI ... 49	CUBTRI ... 453	DLLSQ ... 329
ARTNQ ... 7	CGAMMA ... 61	CURVD ... 413	DLNREL ... 24
ASSGN ... 359	CHEBY ... 375	CURVI ... 413	DMADD ... 179
BADD ... 239	CIN ... 57	CURV1 ... 409	DMCOPY ... 171
BESI ... 79	CIRCV ... 109	CURV2 ... 411	DMCVFS ... 163
BESJ ... 75	CK ... 87	CVBC ... 233	DMCVSF ... 163
BETALN ... 67	CKE ... 88	CVBR ... 233	DMEVAL ... 443
BI ... 85	CKPROD ... 191	CVCB ... 233	DMEXP ... 227
BIE ... 85	CLE ... 229	CVCB1 ... 234	DMFIT ... 441
BLKORD ... 287	CLI ... 59	CVCS ... 255	DMPROD ... 183
BLSQ ... 345	CL1 ... 321	CVDS ... 255	DMSLV ... 195
BPOSE ... 237	CMADD ... 180	CVPRD ... 283	DMSUBT ... 181
BPROD ... 243	CMADJ ... 177	CVPRD1 ... 283	DMTMS ... 183
BRATIO ... 73	CMCONJ ... 173	CVRB ... 233	DNRM2 ... 159
BRCOMP ... 74	CMCOPY ... 171	CVRB1 ... 234	DPADD ... 115
BSLV ... 249	CMCVBS ... 235	CVRS ... 255	DPCHOL ... 207
BSLV1 ... 249	CMCVSB ... 235	CVSC ... 255	DPCOPY ... 113
BSL2 ... 429	CMIMAG ... 169	CVSD ... 255	DPDET ... 219
BSPP ... 425	CMPROD ... 183	CVSR ... 255	DPDIV ... 121
BSSLI ... 79	CMREAL ... 169	DABSLV ... 221	DPINV ... 125
BSSLU ... 75	CMSLV ... 211	DAORD ... 6	DPLE ... 229
BSSLK ... 81	CMSLV1 ... 212	DARTNQ ... 7	DPLPWR ... 123
BSSLY ... 77	CMSUBT ... 182	DASUM ... 157	DPMPAR ... 4
BSJRP ... 423	CMTMS ... 183	DAW ... 47	DPMULT ... 119
BSUBT ... 241	CONSTR ... 17	DAXPY ... 155	DPOSE ... 269
BTPRD ... 245	COSQB ... 371	DBETLN ... 67	DPOSE1 ... 269
BTPRD1 ... 246	COSQF ... 372	DCBCRT ... 139	DPSI ... 66
BTSLV ... 201	COSQI ... 371	DCBRT ... 7	DPSUBT ... 117
BVIP ... 437	CPABS ... 7	DCEIG ... 317	DQAGI ... 450
BVIP2 ... 439	CPOSE ... 271	DCEIGV ... 319	DQAGS ... 447
BVPRD ... 245	CPOSE1 ... 271	DCERF ... 37	DQDCRT ... 139
BVPRD1 ... 245	CPSC ... 381	DCERFC ... 38	DQTCRT ... 139
CAI ... 83	CPSI ... 65	DCGAMA ... 62	DRDVAL ... 94
CAPO ... 11	CREC ... 9	DCMSLV ... 213	DREXP ... 21
CAXPY ... 155	CROUT ... 193	DCOPY ... 145	DRFVAL ... 93
CBADD ... 239	CSADD ... 273	DCPABS ... 7	DRJVAL ... 96
CBCRT ... 139	CSCAL ... 153	DCPOLY ... 141	DRLOG ... 24
CBI ... 83	CSCOPY ... 263	DCPSC ... 382	DRLOG1 ... 24
CBPOSE ... 237	CSEVL ... 129	DCPSI ... 66	DROT ... 149
CBPROD ... 243	CSIMAG ... 259	DCREC ... 9	DROTG ... 13
CBRT ... 7	CSINT ... 405	DCSEVL ... 129	DRPOLY ... 141
CBSLV ... 251	CSINT1 ... 405	DCSQRT ... 9	DSADD ... 273
CBSLV1 ... 251	CSINT2 ... 405	DDOT ... 151	DSCAL ... 153
CBSPL ... 397	CSLOOP ... 407	DDSORT ... 6	DSCOPY ... 263
CBSSLU ... 75	CSLV ... 298	DEI ... 54	DSL ... 294
CBSSLK ... 81	CSLVP ... 215	DEIG ... 305	DSMSLV ... 203
CBSUBT ... 241	CSPROD ... 277	DEIGV ... 307	DSORT ... 6
CBTPD ... 247	CSPSLV ... 297	DEI1 ... 55	DSPROD ... 277
CBTPD1 ... 248	CSREAL ... 259	DELLPI ... 93	DSPSLV ... 293
CBVPD ... 247	CSROT ... 149	DEPI ... 96	DSSUBT ... 275
CBVPD1 ... 247	CSSCAL ... 153	DERF ... 38	DSVDC ... 217
CCOPY ... 145	CSSUBT ... 275	DERFC ... 38	DSVPRD ... 185

DSWAP	...	147	LGRNGV	...	131	PLCCOPY	...	113	SEPDE	...	481
DTASLV	...	223	LGRNGX	...	131	PLEM	...	104	SEVAL	...	403
DTIP	...	175	LLSQ	...	323	PLEM1	...	105	SEVAL1	...	404
DTOPLX	...	209	LLSQMP	...	327	PLPWR	...	123	SEVAL2	...	404
DTPOSE	...	175	LMDIFF	...	353	PMULT	...	119	SFODE	...	473
DTPRD	...	281	LOCPT	...	25	PNDF	...	43	SFODE1	...	476
DTPRD1	...	281	LOPCMP	...	408	PNINV	...	45	SHELL	...	5
DTSLV	...	294	LOPDF	...	408	POCA	...	11	SHELL2	...	5
DVPRD	...	281	LSEI	...	333	PPVAL	...	427	SI	...	57
DVPRD1	...	281	LTRP	...	387	PSI	...	65	SINQB	...	372
EIG	...	301	L2SLV	...	341	PSUBT	...	117	SINQF	...	372
EIGV	...	303	MADD	...	179	QAGI	...	449	SLVMP	...	197
EIGV1	...	303	MCCOPY	...	171	QAGS	...	445	SMADD	...	179
EIG1	...	301	MCVBS	...	235	QDCRT	...	139	SMCOPY	...	171
ELLPF	...	99	MCVDR	...	165	QSUBA	...	446	SMPLX	...	355
ELLPI	...	91	MCVFS	...	163	QTCRT	...	139	SMPROD	...	189
ELPFC1	...	100	MCVRC	...	167	QURV1	...	419	SMSLV	...	203
EPI	...	95	MCVRD	...	165	QURV2	...	420	SMSUBT	...	181
ERF	...	36	MCVSB	...	235	RBND	...	143	SNHCSH	...	19
ERFC	...	36	MCVSF	...	163	RCOMP	...	69	SNRM2	...	159
ERFC1	...	36	MEVAL	...	443	RDVAL	...	92	SPFIT	...	399
ERFINV	...	41	MEXP	...	227	REXP	...	21	SPLIFT	...	398
EXPLI	...	54	MFFT	...	369	RFVAL	...	92	SPLSC	...	347
FFT	...	367	MFFT1	...	369	RISORT	...	5	SPMPAR	...	3
FFT1	...	367	MFIT	...	441	RJVAL	...	95	SPORD	...	285
FMIN	...	349	MKP	...	361	RK	...	477	SPROD	...	277
FRNL	...	51	MPLNMV	...	127	RKF45	...	465	SPSLV	...	289
GAMINV	...	71	MPROD	...	183	RK8	...	479	SQUINT	...	225
GAMLN	...	62	MSLV	...	195	RLOG	...	23	SROT	...	149
GAMMA	...	61	MSUBT	...	181	RLOG1	...	23	SROTG	...	13
GERK	...	469	MTMS	...	183	ROTA	...	11	SSCAL	...	153
GRATIO	...	69	MTPRD	...	279	ROT3	...	15	SSPLX	...	355
HBRD	...	137	MTPRD1	...	279	RPOSE	...	267	SSUBT	...	275
HC	...	31	MVPRD	...	279	RPOSE1	...	267	SSVDC	...	217
HFTI	...	324	MVPRD1	...	279	RSCOPY	...	263	SSWAP	...	147
HFTI2	...	324	NPIVOT	...	194	RSLV	...	290	STLSQ	...	347
HTRP	...	389	NRNG	...	487	SADD	...	273	SURF	...	433
HULL	...	27	NSURF2	...	435	SASUM	...	157	SURF2	...	435
ICAMAX	...	161	ODE	...	461	SAXPY	...	155	SVPRD	...	185
IDAMAX	...	161	OPTF	...	351	SCASUM	...	157	TASLV	...	223
IESLV	...	455	ORTHOS	...	133	SCNRM2	...	159	TIP	...	175
IPMPAR	...	4	ORTHOV	...	133	SCOMP	...	401	TMPROD	...	187
ISAMAX	...	161	ORTHOX	...	134	SCOMP1	...	402	TOPLX	...	209
ISHELL	...	5	PADD	...	115	SCOMP2	...	402	TPOSE	...	175
ISUBX	...	74	PAREA	...	29	SCONJ	...	265	TRP	...	385
KPROD	...	191	PCHOL	...	207	SCOPY	...	145	TSLV	...	290
KROUT	...	194	PCOEFF	...	391	SCVDR	...	257	URNG	...	485
KURVP1	...	417	PDIV	...	121	SCVRC	...	261	VALR2	...	107
KURVP2	...	417	PEQ	...	103	SCVRD	...	257	WNNLS	...	337
KURV1	...	415	PEQ1	...	104	SDOT	...	151	WPFIT	...	395
KURV2	...	416	PFIT	...	393	SEIG	...	309	ZEROIN	...	135
LAINV	...	363	PINV	...	125	SEIGV	...	311			
LE	...	229	PKILL	...	111	SEIGV1	...	311			
LGRNGN	...	131	PKILL3	...	111	SEIG1	...	309			



## DISTRIBUTION

	<u>Copies</u>		<u>Copies</u>
M. Abdulwahab, MS-1M Polaroid Corporation 2 Osborn Street Cambridge, MA 02139	1	Vivien Bowman Technical Information Department Unilever Australia Ltd. P.O. Box 9 Balmain 2041 Reynolds Street Balmain NSW, Australia	1
Dr. Donald C. S. Allison Department of Computer Science Virginia Polytechnic Institute and State University Blacksburg, VA 24061	1	Angela M. Brusca, MS 1297 Westinghouse Defense and Electronics Center P.O. Box 746 Baltimore, MD 21203	1
Donald E. Amos, 5642 Sandia National Laboratories P.O. Box 5800 Albuquerque, NM 87185	1	C. M. Callahan, Code 760 Naval Coastal Systems Laboratory Panama City, FL 32401	1
Rosario Arroyo E.T.S. Ingenieros Industriales Dept. de Energia Nuclear P. de la Castellana, 80 Madrid 6, Spain	1	Mark Chatterton General Mills, Inc. James Ford Bell Technical Center 9000 Plymouth Avenue, North Minneapolis, MN 55427	1
Terence J. Bazou ORI, Inc. 1375 Piccard Drive Rockville, MD 20850	1	H.S. Chen U.S. Army Engineer Waterways Experiment Station Coastal Engineering Research Center Vicksburg, MS 39180	1
Maryann Becker 160/4343 Aerojet ElectroSystems P.O. Box 296 Azusa, CA 91702	1	David Cohoon Damaskos, Inc. Suite 207, Bldg. 200 Chadds Ford Prof. Center Chadds Ford, PA 19317	1
John R. Berg System Development Corporation Satellite Programs Branch 2500 Colorado Avenue Santa Monica, CA 90406	1	Chuck Converse ITT Federal Electric Corporation DS 450 P.O. Box 1886 Vandenberg Air Force Base Lompox, CA 93437	1
Terry D. Boldt IBM Federal Systems Division 2625 Townsgate Rd. Westlake Village, CA 91361	1	Joyce Elliott, Code 4210 Naval Coastal Systems Center Panama City, FL 32407	1
William H. Bonville Federal Electric Corp. Western Space and Missile Center Div. P.O. Box 5728 - Mail Code TE 520 Vandenberg Air Force Base, CA 93437	1	Gerald J. Elias, Code 714 Naval Underwater Systems Center Newport, RI 02840	1

	<u>Copies</u>		<u>Copies</u>
Jules Enig Enig Associates Suite 500 11120 New Hampshire Ave. Silver Spring, MD 20904	1	T. Hengstermann Universitat Oldenburg Postfach 2503 D-2900 Oldenburg, West Germany	1
Brian Feather, MS 450 Calspan Corporation/AEDC Division Arnold Air Force Station, TN 37389	1	Barry E. Herchenroder Department of the Army Coastal Engineering Research Center Kingman Bldg. Fort Belvoir, VA 22060	1
Kirby W. Fong, L-560 Lawrence Livermore Laboratory P.O. Box 5509 Livermore, CA 94550	1	Prof. Roland Horne Sanjay Bhatnagar Dept. of Petroleum Engineering Mitchell Bldg. Stanford University Stanford, CA 94305	1 1
Prof. John E. Francis School of Aerospace and Mech. Engineering University of Oklahoma 865 Asp Ave., Room 212 Norman, OK 73019	1	John L. Houser, MS 2-23 Litton Systems, Inc. Amecon Division 5115 Calvert Road College Park, MD 20740	1
Tom Galib Naval Underwater Systems Center New London, CT 06320	1	Alex Hsia Norton Company Chemical Process Products Division P.O. Box 350 Akron, OH 44309	1
J. D. Gaski Network Analysis Associates, Inc. P.O. Box 8007 Fountain Valley, CA 92728	1	Joshua C. Hung Westinghouse Defense and Electronics Center P.O. Box 746 Baltimore, MD 21203	1
Martin A. Goldberg Webb Institute of Naval Architecture Crescent Beach Road Glen Cove, NY 11542	1	Jan van Kats Rijksuniversiteit Utrecht Academisch Computer Centrum Utrecht Budapestlaan 6 Postbus 80.011 3508 TA Utrecht, Netherlands	1
Donna Hamilton Computing and Communications Services Queen's University Kingston K7L 3N6, Ontario, Canada	1	Ted E. Keller Computer Services Diamond Shamrock Corporation P.O. Box 348 Painesville, OH 44077	1
Richard Hanson IMSL 2500 Park West Tower One Houston, TX 77042-3020	1		
Cheryl Healy, 1-B Polaroid Corporation 750 Main Street Cambridge, MA 02139	1		

	<u>Copies</u>		<u>Copies</u>
Manfred Kory Systems Engineering and Applications Division TRW Defense Systems Group 7600 Colshire Drive McLean, VA 22102	1	Gerald F. Mackey Lee Gantt GTRI/MCSF Electronics Research Bldg. Georgia Institute of Technology Atlanta, GA 30332	1
Eric Krenz Motorola Inc. 1301 E. Algonquin Rd. Schaumburg, IL 60196	1	Eugene Maguin 234 S. Magnolia Lansing, MI 48912	1
Dr. Robbin B. Lake School of Medicine Case Western Reserve University Cleveland, OH 44106	1	Eugene J. Mallove Massachusetts Institute of Technology Lincoln Laboratory P.O. Box 73 Lexington, MA 02173	1
Philipp Lange Bundesanstalt fur Wasserbau Wedeler Landstrasse 157 2000 Hamburg 56 West Germany	1	John M. McKinley Mayo Foundation Medical Sciences Computer Facility 200 First Street, SW Rochester, MN 55905	1
Charles Lawson, MS 301-490 Jet Propulsion Laboratory 4800 Oak Grove Drive Pasadena, CA 91109	1	Chuck Messinger, Code 9122 Naval Ocean Systems Center San Diego, CA 92152	1
Kim Leonard Computer Centre Acadia University Wolfville, Nova Scotia, Canada BOP 1X0	1	Alan Miller CSIRO Division of Mathematics and Statistics Bayview Ave. Clayton, Victoria, Australia 3168	1
Dr. Stephen Lipscomb Mathematics Department Mary Washington College Fredericksburg, VA 22401	1	Keith Mize Dept. of Physics and Astronomy University of South Carolina Columbia, SC 29208	1
Jim Liverman Naval Weapons Station NQEC - Code 302 Yorktown, VA 23691	1	Dr. Richard E. Nance Department of Computer Science Virginia Polytechnic Institute and State University Blacksburg, VA 24061	1
Dr. Luigs Internationale Atomreaktorbau GmbH 9130 EDV und Mathematik Postfach 5060 Bergisch Gladbach 1, West Germany	1	K. D. Needham Easams LTD Lyon Way, Frimley Road Camberley GU16 5EX, Surrey, England	1

	<u>Copies</u>		<u>Copies</u>
Steve Orbon Westinghouse Electric Corporation Box 158 Madison, PA 15663	1	C.K. Sollitt, Director Hinsdale Wave Research Laboratory Oregon State University Corvallis, OR 97331	1
Alvin B. Owens, Code 2303 Naval Research Laboratory Washington, DC 20375	1	Dr. James F. Steffe Dept. of Agricultural Engineering Michigan State University East Lansing, MI 48824	1
George Perez Research Computational Division Naval Research Laboratory Washington, DC 20375	1	Fred Stern Science Applications, Inc. 134 Holiday Court, Suite 318 Annapolis, MD 21401	1
Scott M. Ramsey, Code EB2C Pratt and Whitney Aircraft 400 Main Street East Hartford, CT 06108	1	Connie A. Stirgus Department of the Army Experiment Station, Corps of Engineers P.O. Box 631 Vicksburg, MS 39180	1
John Reed, Code A1 16441 Space Center Blvd. Houston, TX 77058	1	Vish Subramaniam Westinghouse Electric Corporation Avenue A and West Street, Forrest Hills Pittsburg, PA 15221	1
Steven Reuman Polaroid Corp. 750 Main St., 3rd Floor Cambridge, MA 02139	1	Dan Tarrant, Code 8511 Naval Air Development Center Warminster, PA 18974	1
Vitor Rodrigues Correios e Telecomunicacoes de Portugal R.S. Jose, 10-70 1198 Lisboa Codex, Portugal	1	Neil Toews Energy, Mines, and Resources Canada Mining Research Laboratories 555 Booth Street Ottawa, Ontario K1A0G1, Canada	1
Peter Rogers Mechanical Engineering, SSTC Georgia Institute of Technology Atlanta, GA 30332	1	D.S. Toomb Aerojet ElectroSystems Company 1100 West Hollyvale Street P.O. Box 296-III Azusa, CA 91702	1
Matthew R. Scalzetti General Electric Co. Radar Systems Dept. P.O. Box 4840 Syracuse, NY 13221	1	T.G. Truong Technical Systems, Vogel Computing Services Ministry of Works and Development P.O. Box 12-041 Wellington, New Zealand	1
Ing. H.P. Schwefel Universitat Dortmund FB Informatik, Systemanalyse Postfach 500500 4600 Dortmund 50, West Germany	1		

	<u>Copies</u>		<u>Copies</u>
Dr. Carlton W. Ulbrich Dept. of Physics and Astronomy Clemson University Clemson, SC 29634	1	Commander, Naval Sea Systems Command Attn: SEA 05R Department of the Navy Washington, DC 20362	1
James H. Verner Dept. of Mathematics and Statistics Queen's University Kingston K7L 3N6, Ontario, Canada	1	U.S. Naval Observatory 34th St. and Massachusetts Ave., N.W. Attn: Library Washington, DC 20390	1
Todd L. Walton Dept. of the Army Experiment Station, Corps of Engineers P.O. Box 631 Vicksburg, MS 39180	1	Director, U.S. Army Ballistic Research Laboratory Attn: DRDAR-TSB-S (Stinfo) Aberdeen Proving Ground, MD 21005	1
Dr. Layne T. Watson Department of Computer Science Virginia Polytechnic Institute and State University Blacksburg, VA 24061	1	Director, National Bureau of Standards Attn: Library Gaithersburg, MD 20760	1
Munib Wober Polaroid Corp. 750 Main St., 3rd Floor Cambridge, MA 02139	1	Library of Congress Attn: Exchange and Gift Division Washington, DC 20540	1
Jan York Norton Company Chemical Process Products Division P.O. Box 350 Akron, OH 44309	1	Internal Distribution:	
		D	1
		D1	1
		D2	1
		D21	1
		D22	1
		D23	1
		D24	1
		D25	1
		D4	1
Chief of Naval Operations		E	1
Attn: OP-961	1	E02	1
OP-982	1	E03	1
OP-983	1	E04	1
OP-987	1	E10	1
Washington, DC 20350		E20	1
		E211	1
Commander, Naval Facilities Engineering Command		E30	1
Attn: Code 032F	1	E40	1
200 Stovall St.		E43	800
Alexandria, VA 22332		E50	1
		F	1
Office of Naval Research	1	F02	1
Code 411		F10	1
800 North Quincy St.		F20	1
Arlington, VA 22217		F30	1
		F40	1

Copies

G	1
G02	1
G10	1
G20	1
G30	1
G40	1
G60	1
H	1
H10	1
H20	1
H30	1
K	1
K02	1
K10	1
K20	1
K40	1
K50	1
N	1
N02	1
N10	1
N20	1
N30	1
N40	1
R	1
R02	1
R10	1
R30	1
R40	1
U	1
U01	1
U02	1
U10	1
U20	1
U30	1
U40	1

REPORT DOCUMENTATION PAGE			Form Approved OMB No. 0704-0188	
Public reporting burden for this collection of information is estimated to average 1 hour per response, including the time for reviewing instructions, searching existing data sources, gathering and maintaining the data needed, and completing and reviewing the collection of information. Send comments regarding this burden estimate or any other aspect of this collection of information, including suggestions for reducing this burden, to Washington Headquarters Services, Directorate for Information Operations and Reports, 1215 Jefferson Davis Highway, Suite 1204, Arlington, VA 22202-4302, and to the Office of Management and Budget, Paperwork Reduction Project (0704-0188), Washington, DC 20503.				
1. AGENCY USE ONLY (Leave blank)	2. REPORT DATE January 1990	3. REPORT TYPE AND DATES COVERED		
4. TITLE AND SUBTITLE NSWC LIBRARY OF MATHEMATICS SUBROUTINES		5. FUNDING NUMBERS		
6. AUTHOR(S) Alfred H. Morris, Jr.				
7. PERFORMING ORGANIZATION NAME(S) AND ADDRESS(ES) Naval Surface Warfare Center (E43) Dahlgren, VA 22448-5000		8. PERFORMING ORGANIZATION REPORT NUMBER NSWC TR 90-21		
9. SPONSORING/MONITORING AGENCY NAME(S) AND ADDRESS(ES)		10. SPONSORING/MONITORING AGENCY REPORT NUMBER		
11. SUPPLEMENTARY NOTES				
12a. DISTRIBUTION/AVAILABILITY STATEMENT Approved for public release; distribution unlimited.		12b. DISTRIBUTION CODE		
13. ABSTRACT (Maximum 200 words)  The NSWC library is a library of general-purpose Fortran subroutines that provide a basic computational capability in a variety of mathematical activities. Emphasis has been placed on the transportability of the codes. Subroutines are available in the following areas: Elementary Operations, Geometry, Special Functions, Polynomials, Vectors, Matrices, Large Dense Systems of Linear Equations, Banded Matrices, Sparse Matrices, Eigenvalues and Eigenvectors, $\ell_1$ Solution of Linear Equations, Least-Squares Solution of Linear Equations, Optimization, Transforms, Approximation of Functions, Curve Fitting, Surface Fitting, Manifold Fitting, Numerical Integration, Integral Equations, Ordinary Differential Equations, Partial Differential Equations, and Random Number Generation.				
14. SUBJECT TERMS		15. NUMBER OF PAGES		
		16. PRICE CODE		
17. SECURITY CLASSIFICATION OF REPORT UNCLASSIFIED	18. SECURITY CLASSIFICATION OF THIS PAGE UNCLASSIFIED	19. SECURITY CLASSIFICATION OF ABSTRACT UNCLASSIFIED	20. LIMITATION OF ABSTRACT	

## GENERAL INSTRUCTIONS FOR COMPLETING SF 298

The Report Documentation Page (RDP) is used in announcing and cataloging reports. It is important that this information be consistent with the rest of the report, particularly the cover and its title page. Instructions for filling in each block of the form follow. It is important to *stay within the lines* to meet optical scanning requirements.

**Block 1. Agency Use Only (Leave blank).**

**Block 2. Report Date.** Full publication date including day, month, and year, if available (e.g. 1 Jan 88). Must cite at least the year.

**Block 3. Type of Report and Dates Covered.** State whether report is interim, final, etc. If applicable, enter inclusive report dates (e.g. 10 Jun 87 - 30 Jun 88).

**Block 4. Title and Subtitle.** A title is taken from the part of the report that provides the most meaningful and complete information. When a report is prepared in more than one volume, repeat the primary title, add volume number, and include subtitle for the specific volume. On classified documents enter the title classification in parentheses.

**Block 5. Funding Numbers.** To include contract and grant numbers; may include program element number(s), project number(s), task number(s), and work unit number(s). Use the following labels:

C - Contract	PR - Project
G - Grant	TA - Task
PE - Program Element	WU - Work Unit Accession No.

**BLOCK 6. Author(s).** Name(s) of person(s) responsible for writing the report, performing the research, or credited with the content of the report. If editor or compiler, this should follow the name(s).

**Block 7. Performing Organization Name(s) and address(es).** Self-explanatory.

**Block 8. Performing Organization Report Number.** Enter the unique alphanumeric report number(s) assigned by the organization performing the report.

**Block 9. Sponsoring/Monitoring Agency Name(s) and Address(es).** Self-explanatory.

**Block 10. Sponsoring/Monitoring Agency Report Number.** (If Known)

**Block 11. Supplementary Notes.** Enter information not included elsewhere such as: Prepared in cooperation with...; Trans. of...; To be published in... . When a report is revised, include a statement whether the new report supersedes or supplements the older report.

**Block 12a. Distribution/Availability Statement.**

Denotes public availability or limitations. Cite any availability to the public. Enter additional limitations or special markings in all capitals (e.g. NOFORN, REL, ITAR).

DOD - See DoDD 5230.24, "Distribution Statements on Technical Documents."  
DOE - See authorities.  
NASA - See Handbook NHB 2200.2  
NTIS - Leave blank

**Block 12b. Distribution Code.**

DOD - Leave blank.  
DOE - Enter DOE distribution categories from the Standard Distribution for Unclassified Scientific and Technical Reports.  
NASA - Leave blank.  
NTIS - Leave blank.

**Block 13. Abstract.** Include a brief (*Maximum 200 words*) factual summary of the most significant information contained in the report.

**Block 14. Subject Terms.** Keywords or phrases identifying major subjects in the report.

**Block 15. Number of Pages.** Enter the total number of pages.

**Block 16. Price Code.** Enter appropriate price code (*NTIS only*)

**Block 17.-19. Security Classifications.** Self-explanatory. Enter U.S. Security Classification in accordance with U.S. Security Regulations (i.e., UNCLASSIFIED). If form contains classified information, stamp classification on the top and bottom of this page.

**Block 20. Limitation of Abstract.** This block must be completed to assign a limitation to the abstract. Enter either UL (unlimited or SAR (same as report)). An entry in this block is necessary if the abstract is to be limited. If blank, the abstract is assumed to be unlimited.